

Data Science and Machine Learning

Prof. Michalis Vlachos

The Data Science process

Example 1: What are the demographics of our customers?

2

Example 2: What is the sentiment about our company in twitter? Mostly positive or negative?



Regression – Predicting numeric values

3

- Regression is an instance of a **supervised learning** algorithm
- We are given a set of features for some objects/entities (customers, products, etc) AND also the **numeric value** of what we want to predict.

Age	Income	#kids	Withdrawal amount next month
33	100	1	15
55	150	3	21

- Given a new object/entity and its feature values. What would be the numeric value that we want to predict?

38	90	2	??
----	----	---	----

Applications - Examples

4

- Given [square meters, number of bathrooms, ...] \rightarrow ? house price



- Given [size of engine, weight of car, ...] \rightarrow ? km/liter of gas



Discussion:

- Find some applications of regression

Some examples of regression

6

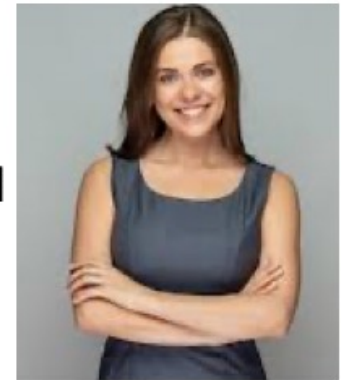
- Regression (we predict a number)



→ 60 ,



→ 42, ... } and then when a new example arrived

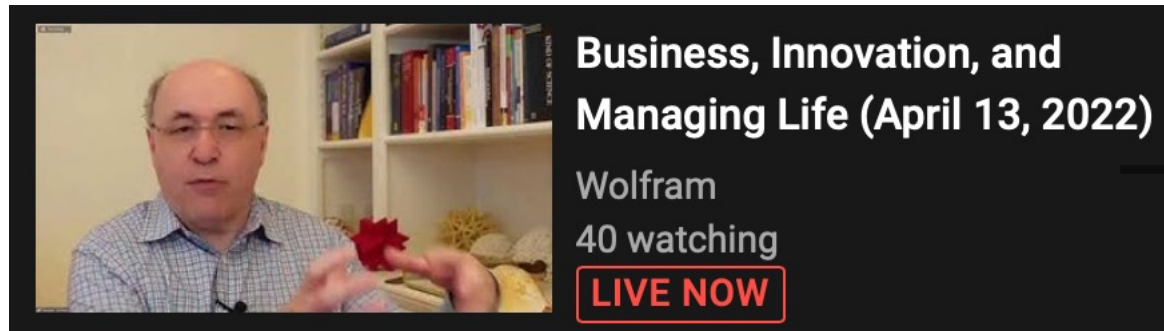


→ ??

Some examples of regression

7

- Regression (we predict a number)



→ IQ = 160



→ IQ = ??

What input features would you use for this application?

Let's recall (some terminology)

Observation

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	OnlineSecurity	...	TotalCharges
0	7590-VHVEG	Female	0	Yes	No	1	No	No	...	29.85
1	5575-GNVDE	Male	0	No	No	34	Yes	Yes	...	1889.5
2	3668-QPYBK	Male	0	No	No	2	Yes	Yes	...	108.15
3	7795-CFOCW	Male	0	No	No	45	No	Yes	...	1840.75
4	9237-HQITU	Female	0	No	No	2		No	...	151.65

Observation could be a customer, a patient, a car, a country, a novel, a drug, a movie etc

Observation = one row

Feature or attribute

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	OnlineSecurity	...	TotalCharges
0	7590-VHVEG	Female	0	Yes	No	1	No	No	...	29.85
1	5575-GNVDE	Male	0	No	No	34	Yes	Yes	...	1889.5
2	3668-QPYBK	Male	0	No	No	2	Yes	Yes	...	108.15
3	7795-CFOCW	Male	0	No	No	45	No	Yes	...	1840.75
4	9237-HQITU	Female	0	No		2	Yes	No	...	151.65

Feature x = column
(independent variables or predictors)

Feature or attribute

customerID		gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	OnlineSecurity	...	TotalCharges
0	7590-VHVEG	Female	0	Yes	No	1	No	No	...	29.85
1	5575-GNVDE	Male	0	No	No	34	Yes	Yes	...	1889.5
2	3668-QPYBK	Male	0	No	No	2	Yes	Yes	...	108.15
3	7795-CFOCW	Male	0	No	No	45	No	Yes	...	1840.75
4	9237-HQITU	Female	0	No		2	Yes	No	...	151.65

All the features → X

Feature or attribute

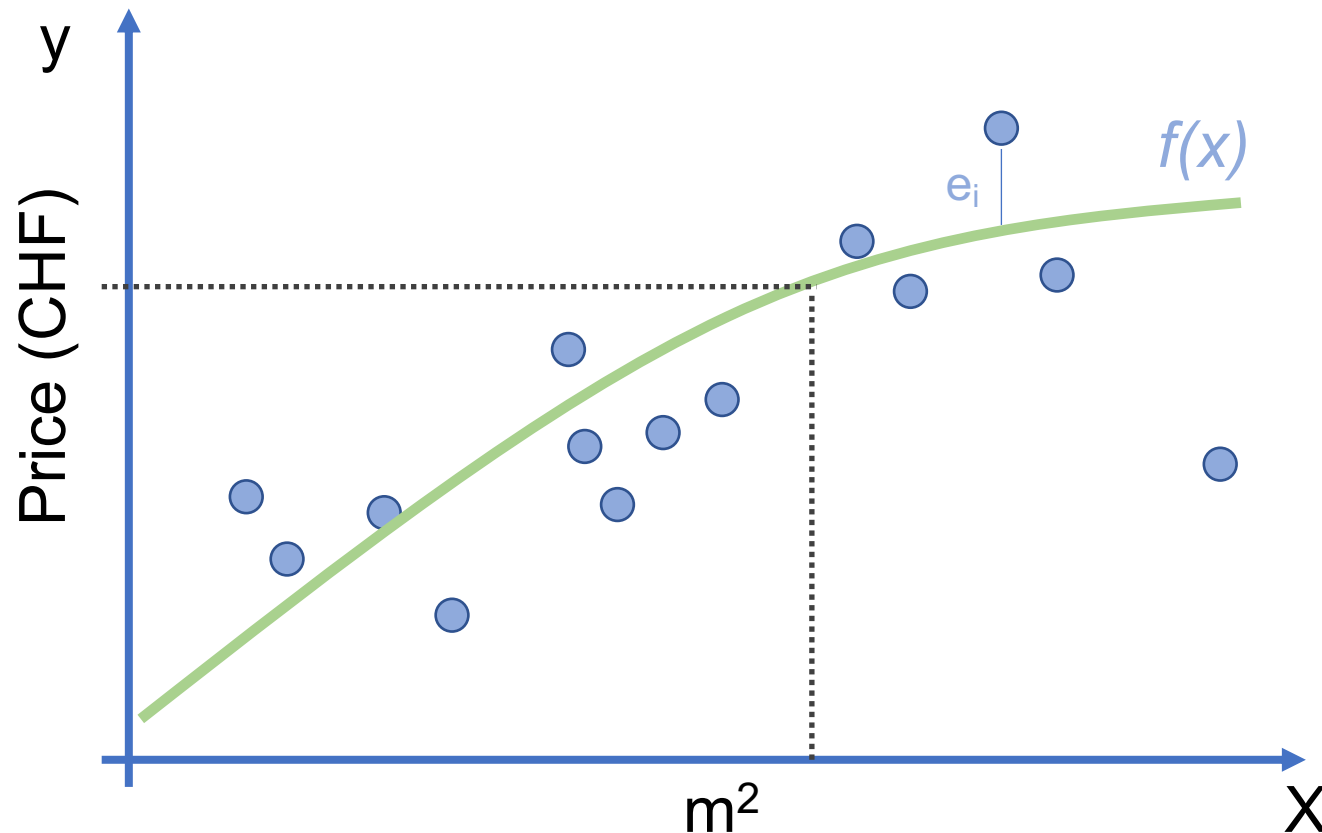
	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	OnlineSecurity	...	TotalCharges
0	7590-VHVEG	Female	0	Yes	No	1	No	No	...	29.85
1	5575-GNVDE	Male	0	No	No	34	Yes	Yes	...	1889.5
2	3668-QPYBK	Male	0	No	No	2	Yes	Yes	...	108.15
3	7795-CFOCW	Male	0	No	No	45	No	Yes	...	1840.75
4	9237-HQITU	Female	0	No	No	2	Yes	No	...	151.65

Target variable y
(dependent variable)

Model representation – Regression

A motivating example

- I would like to sell my house/apartment
- How much should I sell it?

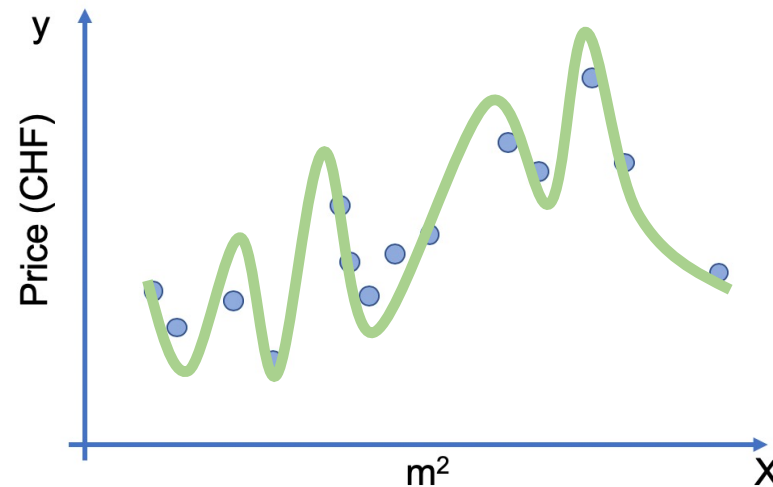
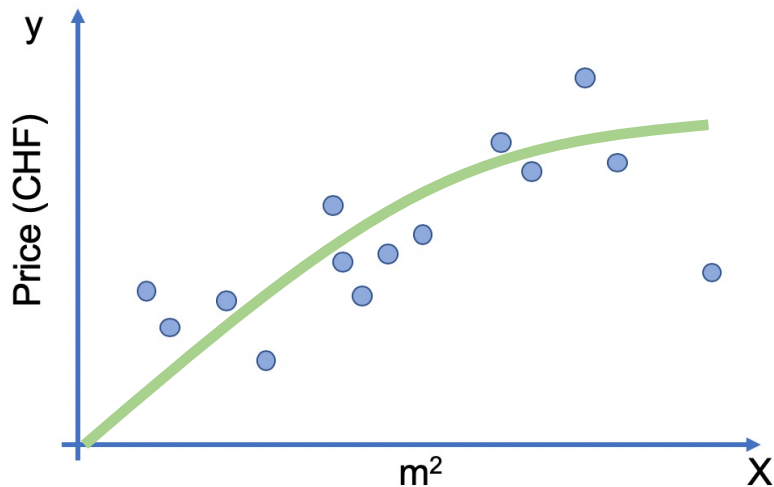
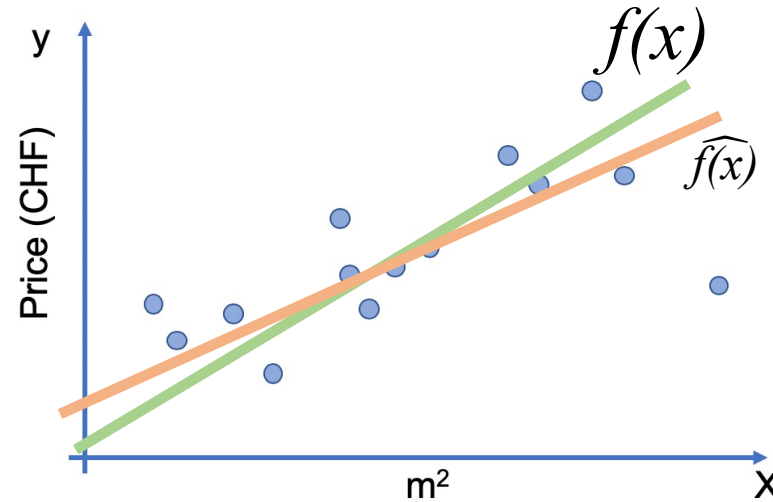
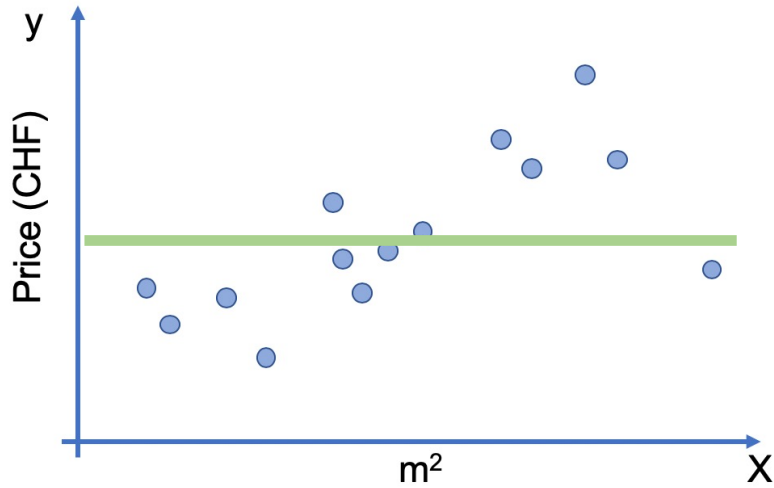


House area	Sale Price
120	930,000
65	705,000
154	2,010,000
...	...

Regression model:
 $y_i = f(x_i) + e_i$

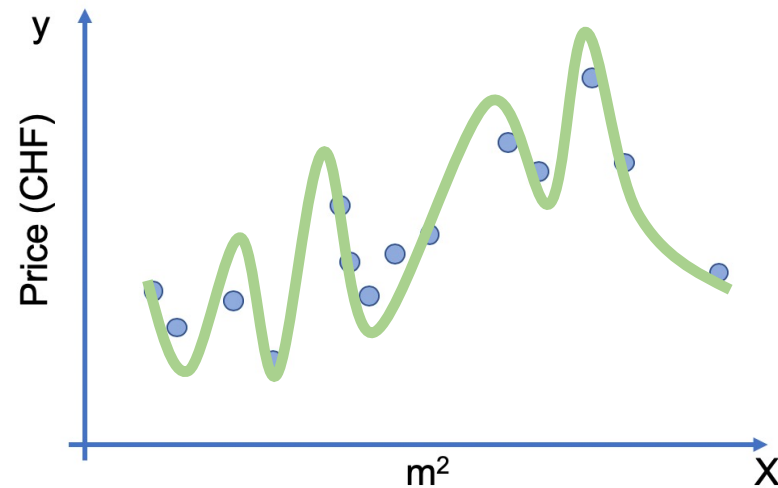
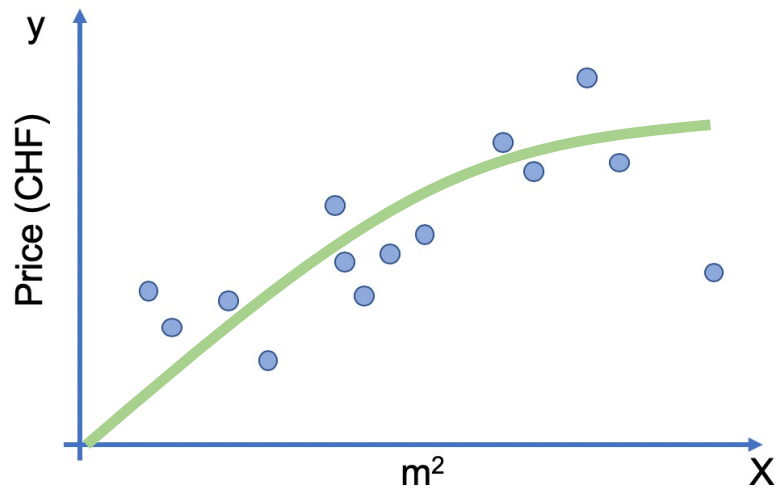
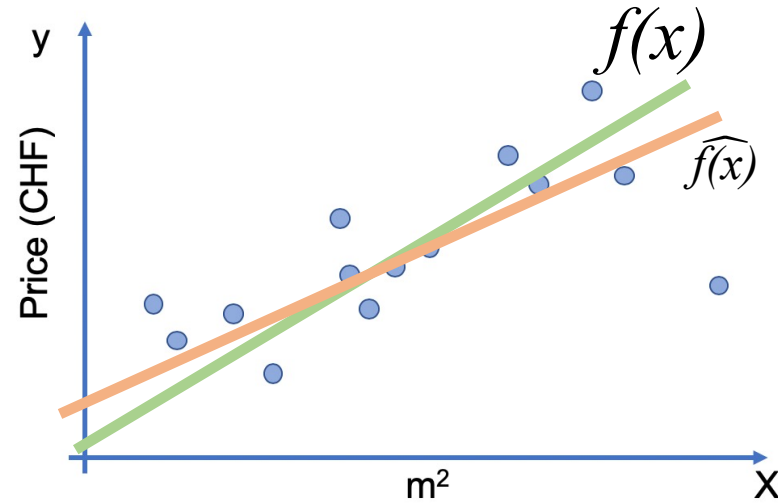
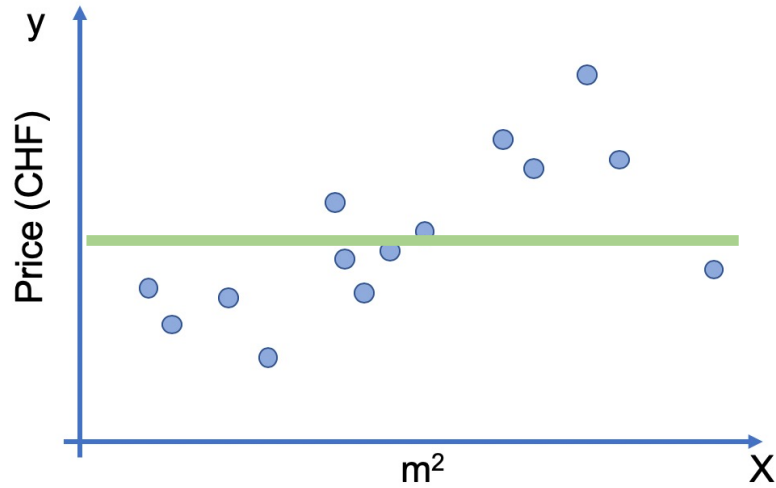
What is a good model $f(x)$?

15

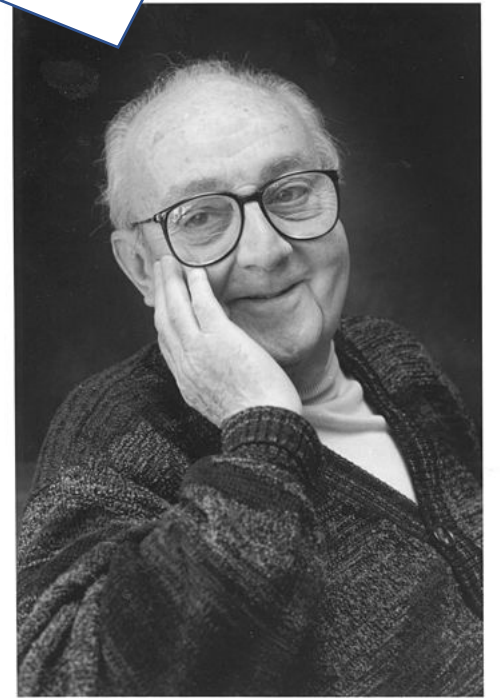


For now we will assume that there exists *linear model* $f(x)$ and we will try to approximate it with an $\hat{f}(x)$ function which we learn from the data.

What is a good model $f(x)$?



All models are
wrong, but some
are useful!



George Box
1919 - 2013

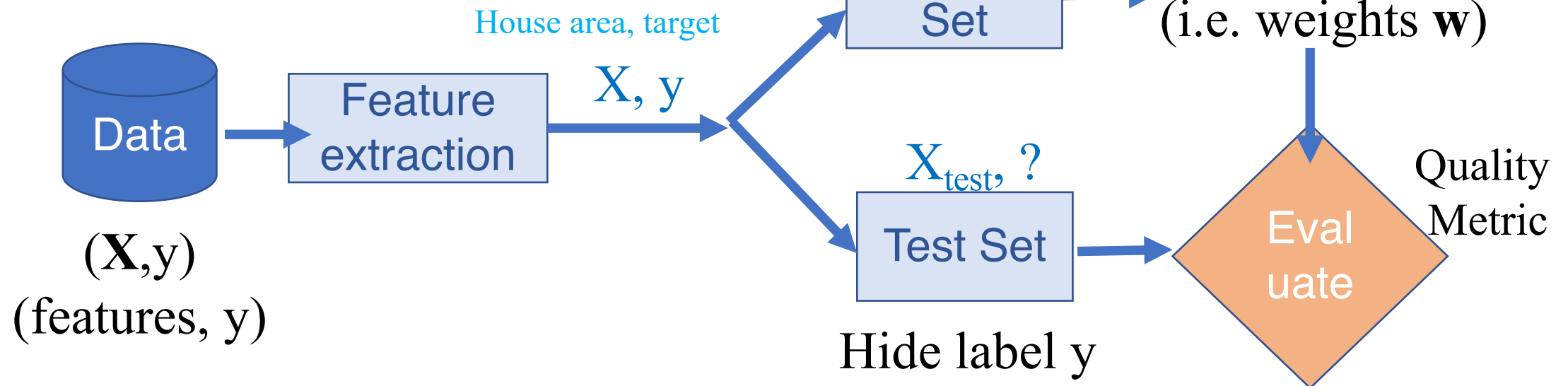
Building the predictive model

17

House area	Lot area	# baths	Sale Price
120	500	...	930k
65	350	...	705k
154	0	...	2010k
...

Data

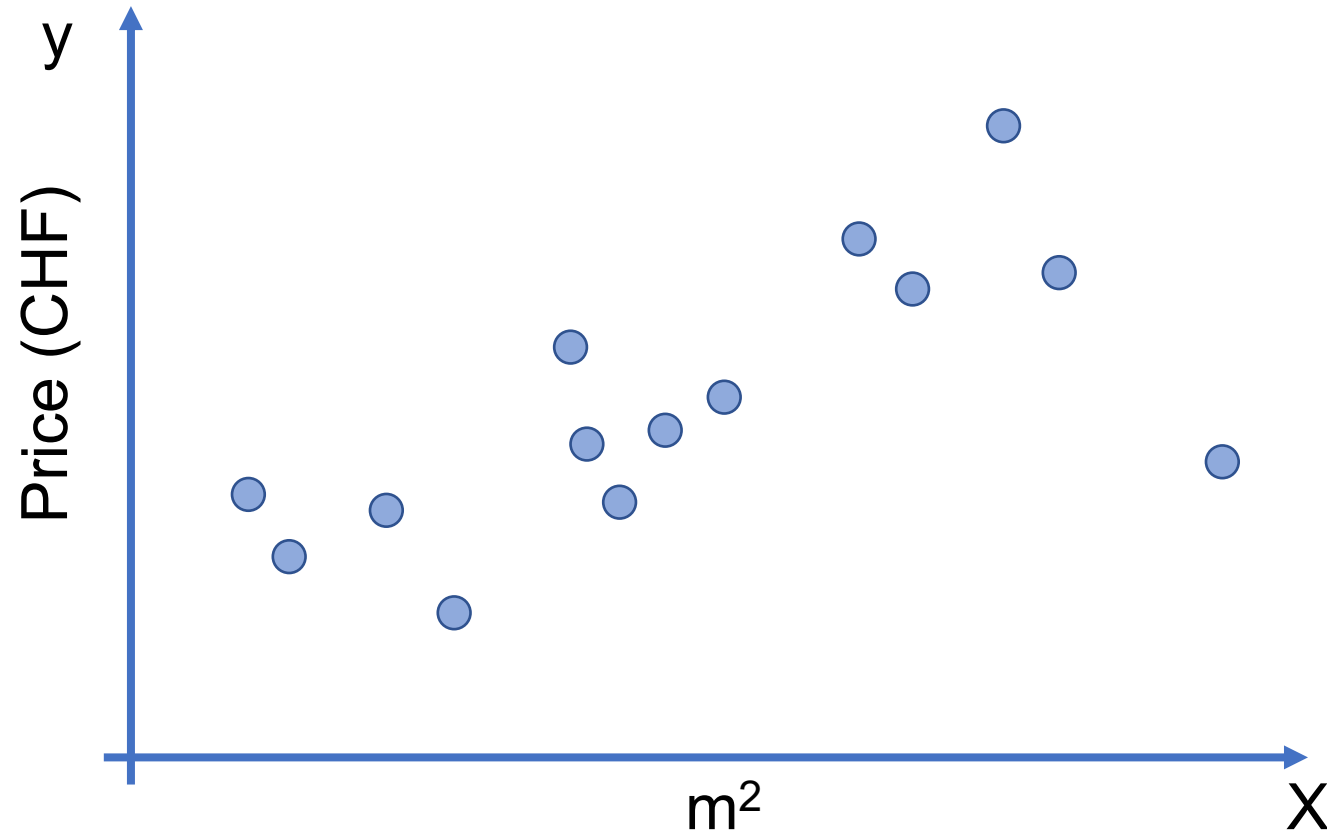
\mathbf{X} y
actual sale



Model representation – Linear Regression

Simple linear regression model

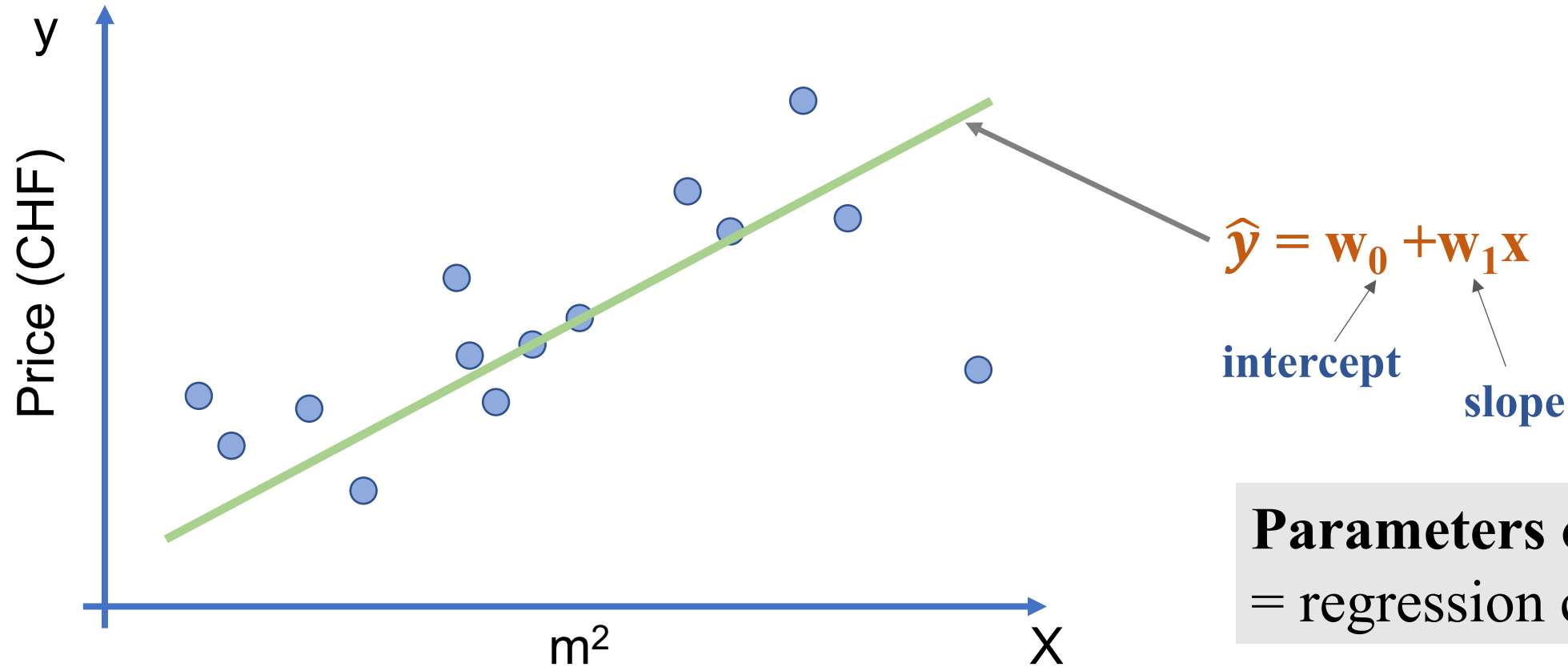
19



Simple linear regression model

20

Fit a linear line through the data

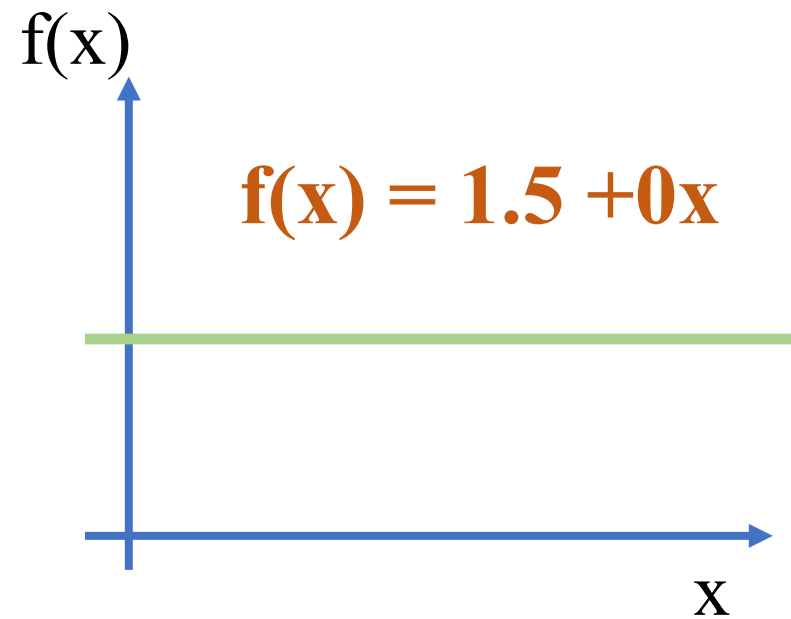


Parameters of the model
= regression coefficients

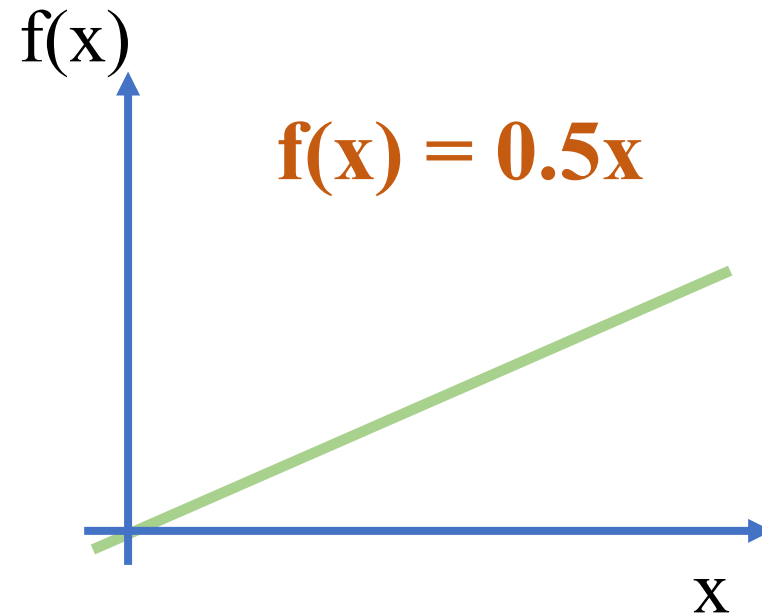
(...let's remember...)

21

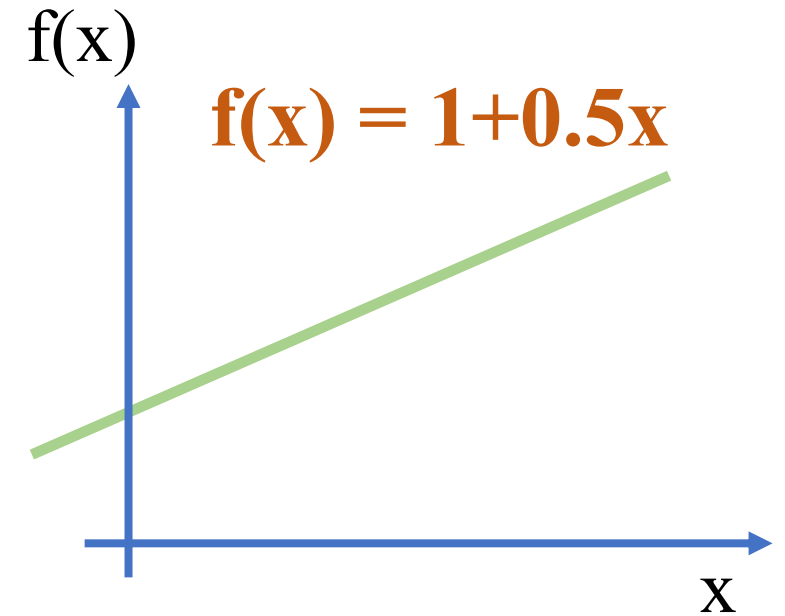
$$f(x) = w_0 + w_1 x$$



$$w_0 = 1.5$$
$$w_1 = 0$$



$$w_0 = 0$$
$$w_1 = 0.5$$



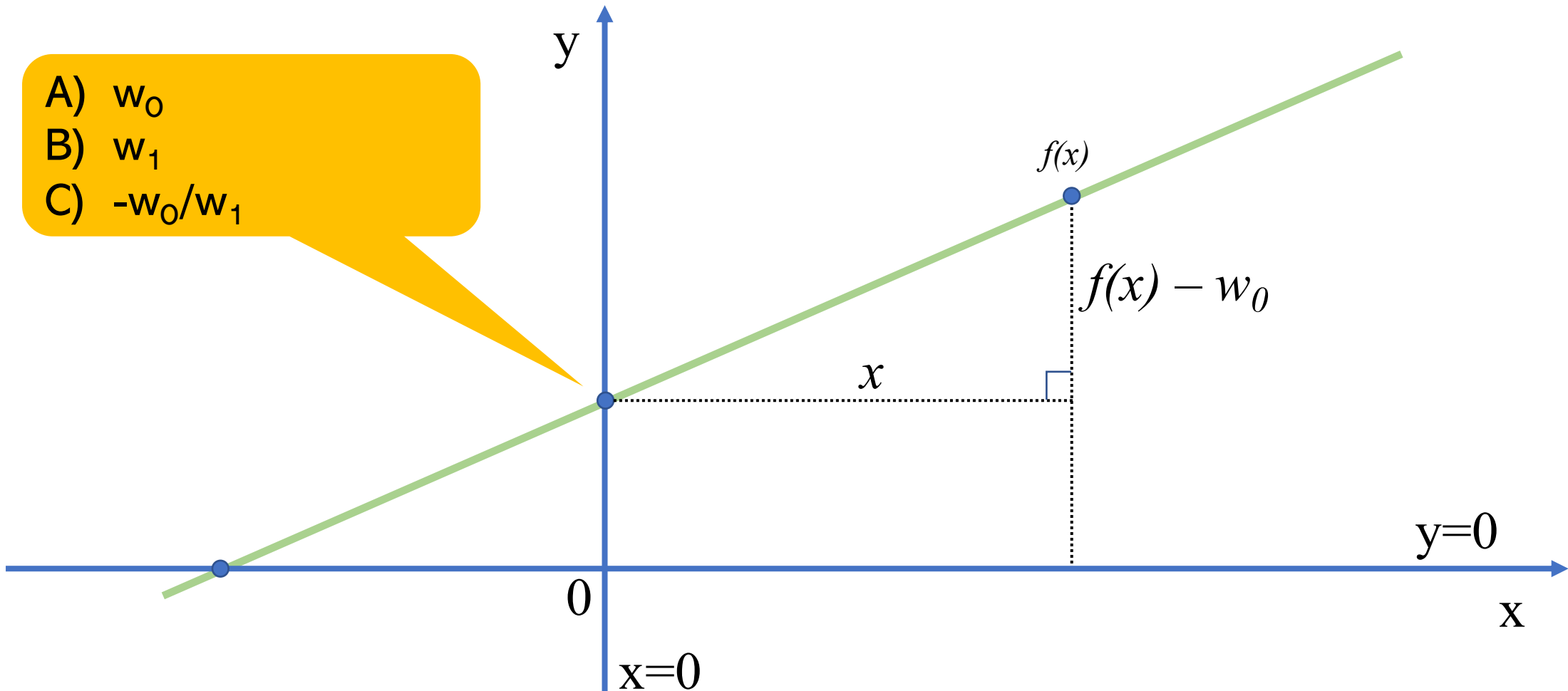
$$w_0 = 1$$
$$w_1 = 0.5$$

(...let's remember...)

22

$$f(x) = w_0 + w_1 x$$

- A) w_0
- B) w_1
- C) $-w_0/w_1$

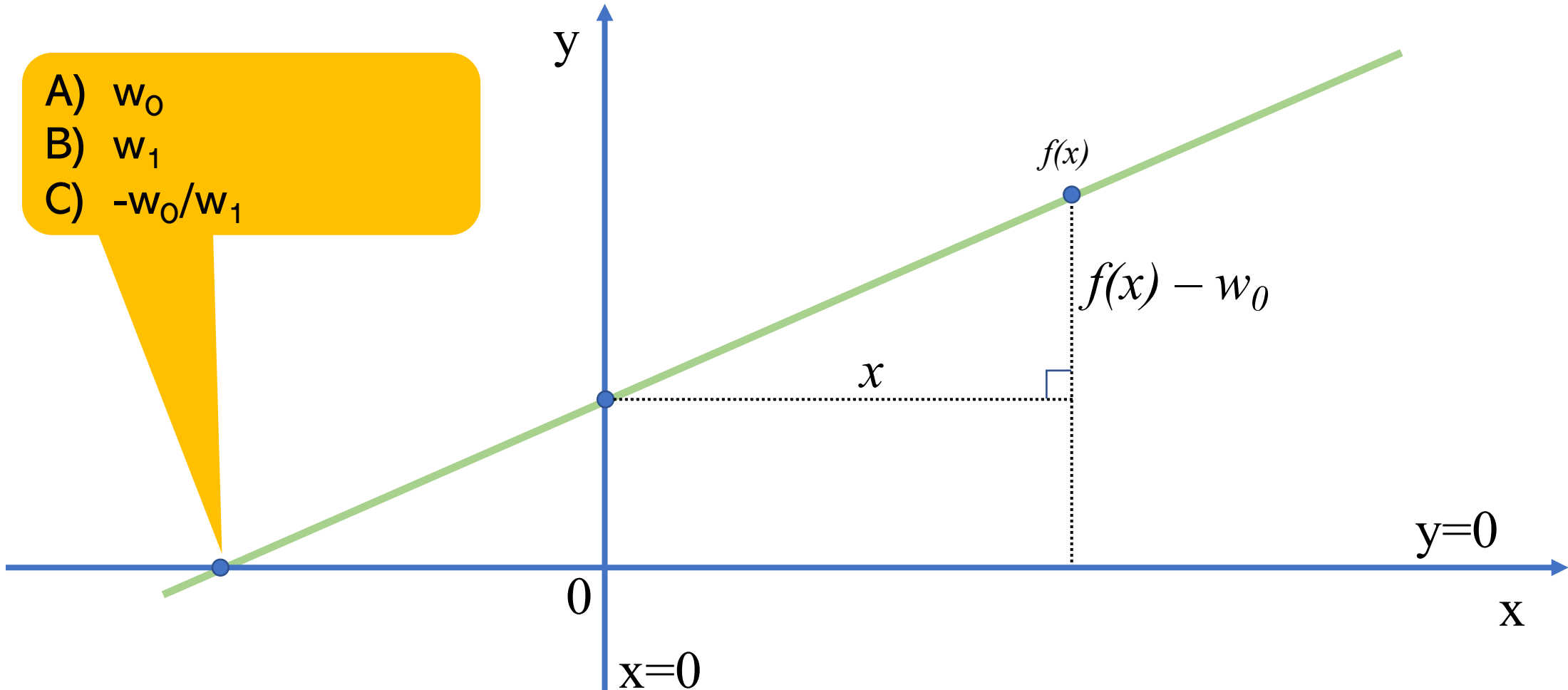


(...let's remember...)

23

$$f(x) = w_0 + w_1 x$$

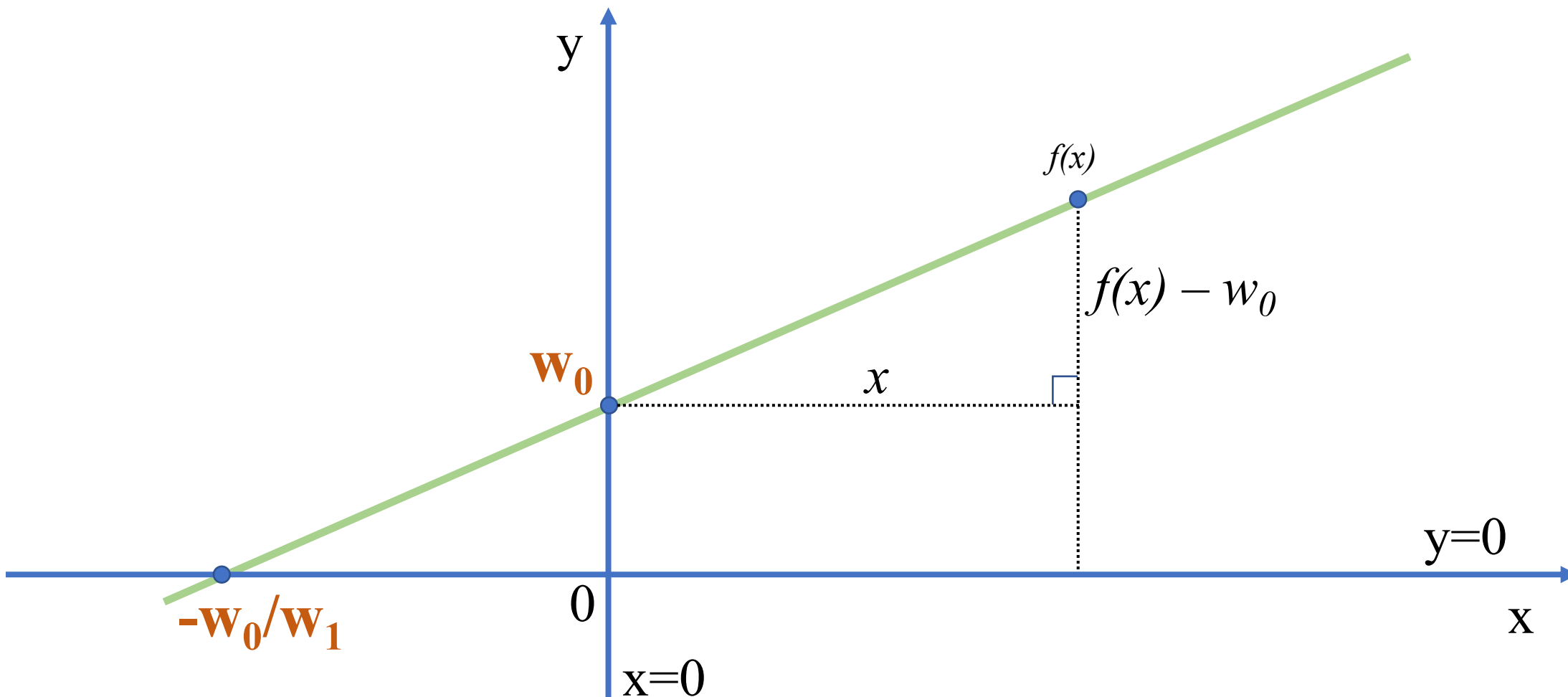
- A) w_0
- B) w_1
- C) $-w_0/w_1$



(...let's remember...)

24

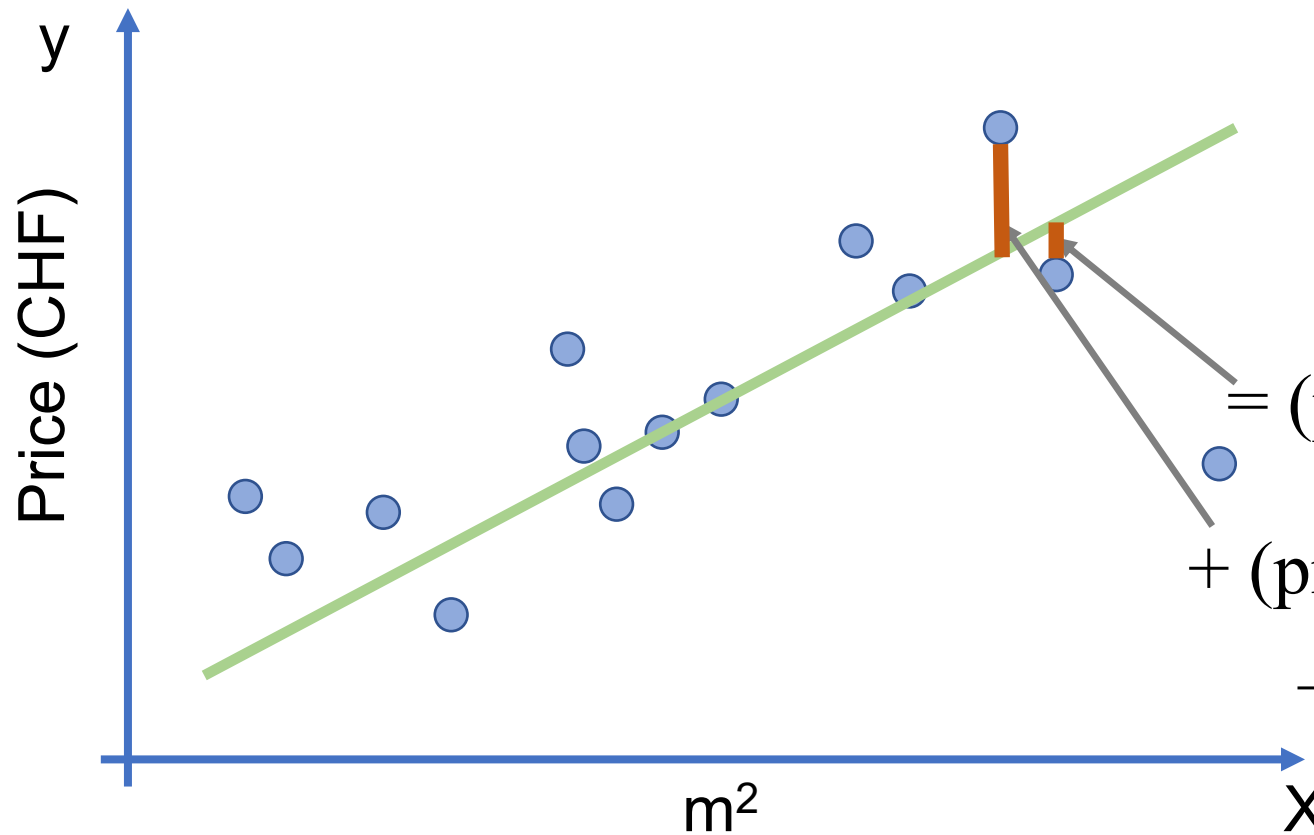
$$f(x) = w_0 + w_1 x$$



Error of using some line model

25

Residual sum of squares (RSS)

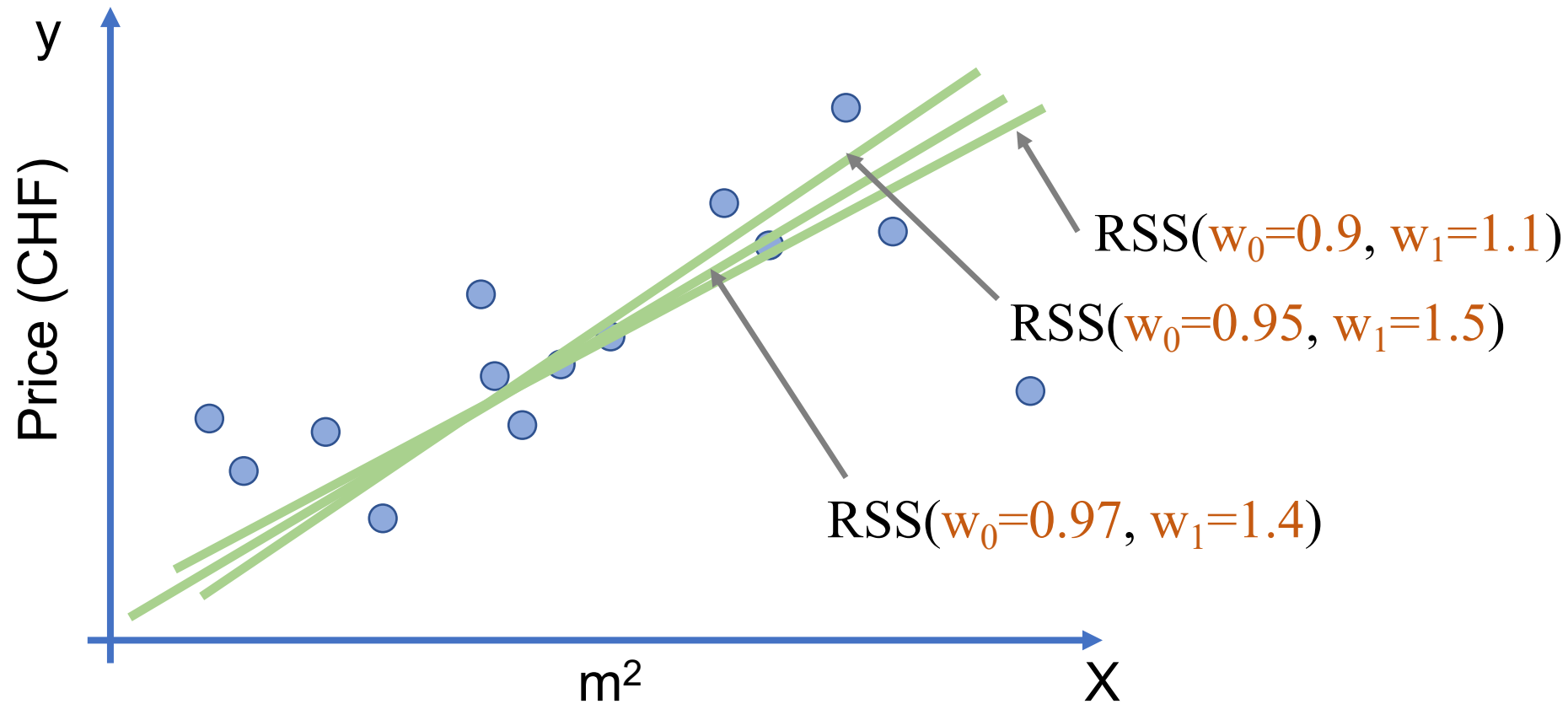


$$\text{RSS}(\mathbf{w}_0, \mathbf{w}_1) = \sum_{i=1}^N (\underbrace{y_i}_{\text{actual value}} - \underbrace{(\mathbf{w}_0 + \mathbf{w}_1 x_i)}_{\text{predicted value}})^2$$
$$= (\text{price_house1} - (\mathbf{w}_0 + \mathbf{w}_1 \text{house_area1}))^2$$
$$+ (\text{price_house_2} - (\mathbf{w}_0 + \mathbf{w}_1 \text{house_area2}))^2$$
$$+ \dots$$

Find the best line

26

Minimize the cost across all possible w_0, w_1

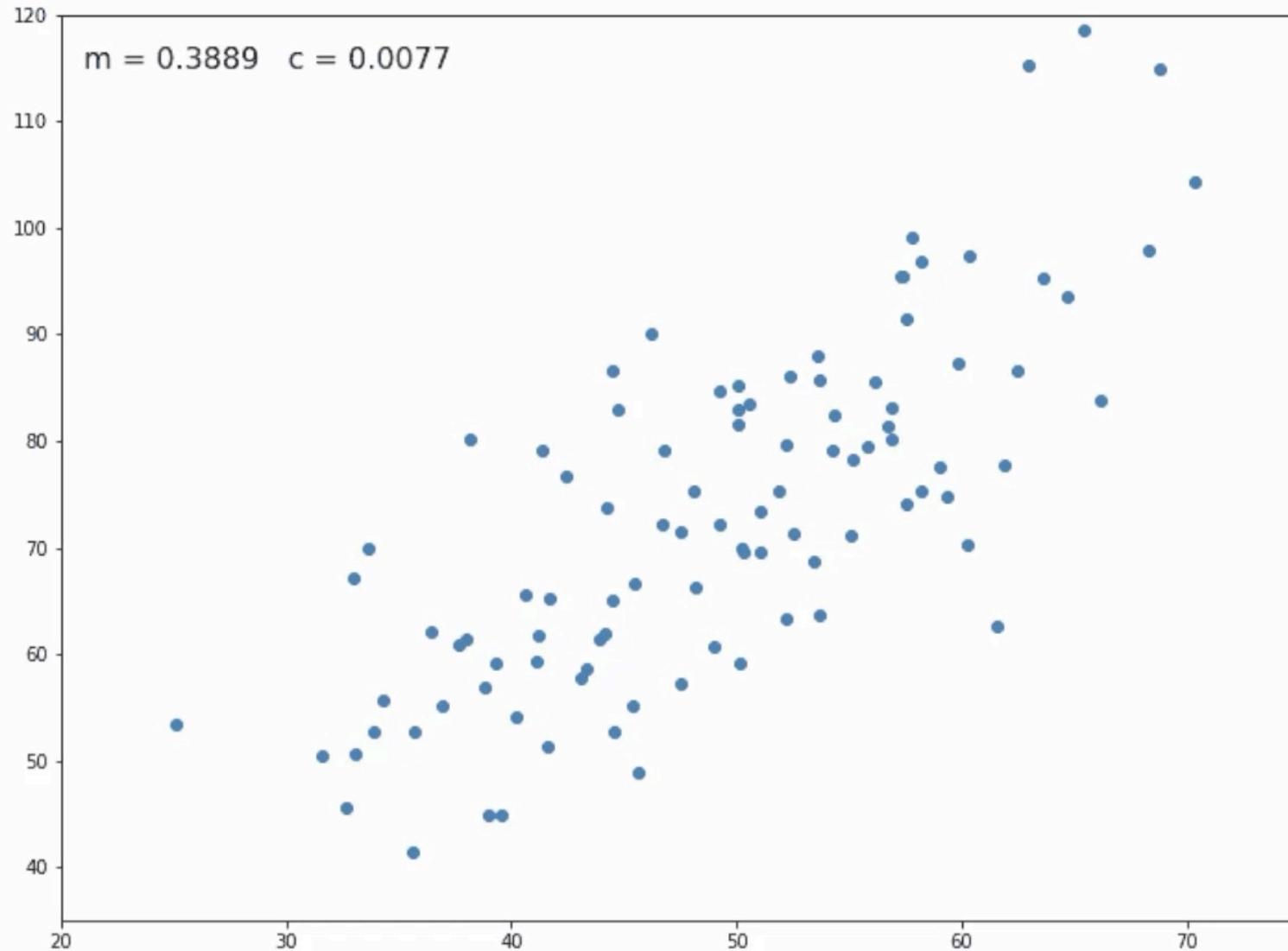


How to find the best line – Gradient Descent

27

- Typically, we find the best line using a method called “**gradient descent**”. This is an iterative method that tries different solutions and guides the tries based on the error (RSS in this case).
- Gradient descent does not try *randomly* for different solutions, but tries values that look most promising and tries in every step to improve the previous estimate.

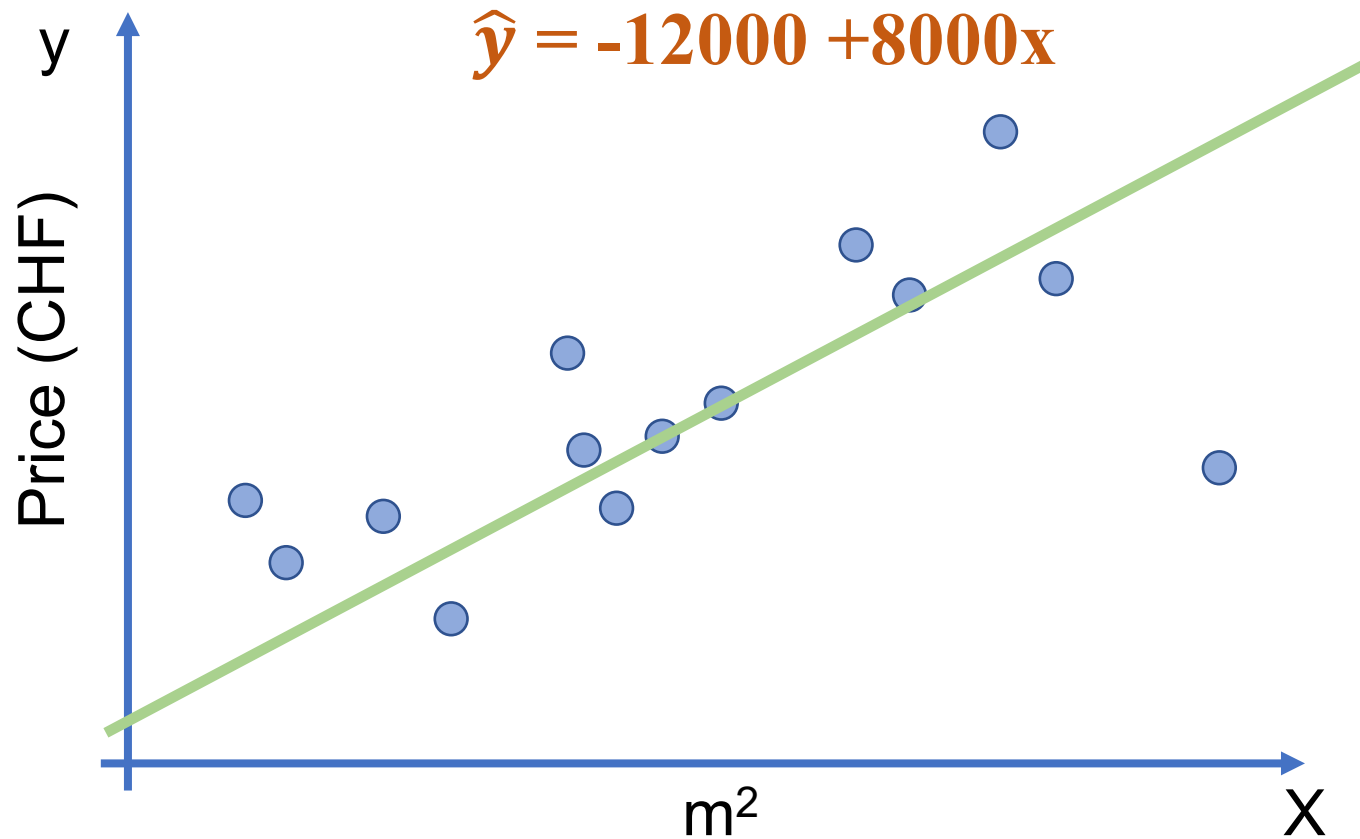
Demo of gradient descent



Doing predictions

29

Assume this is the model with the lowest RSS cost



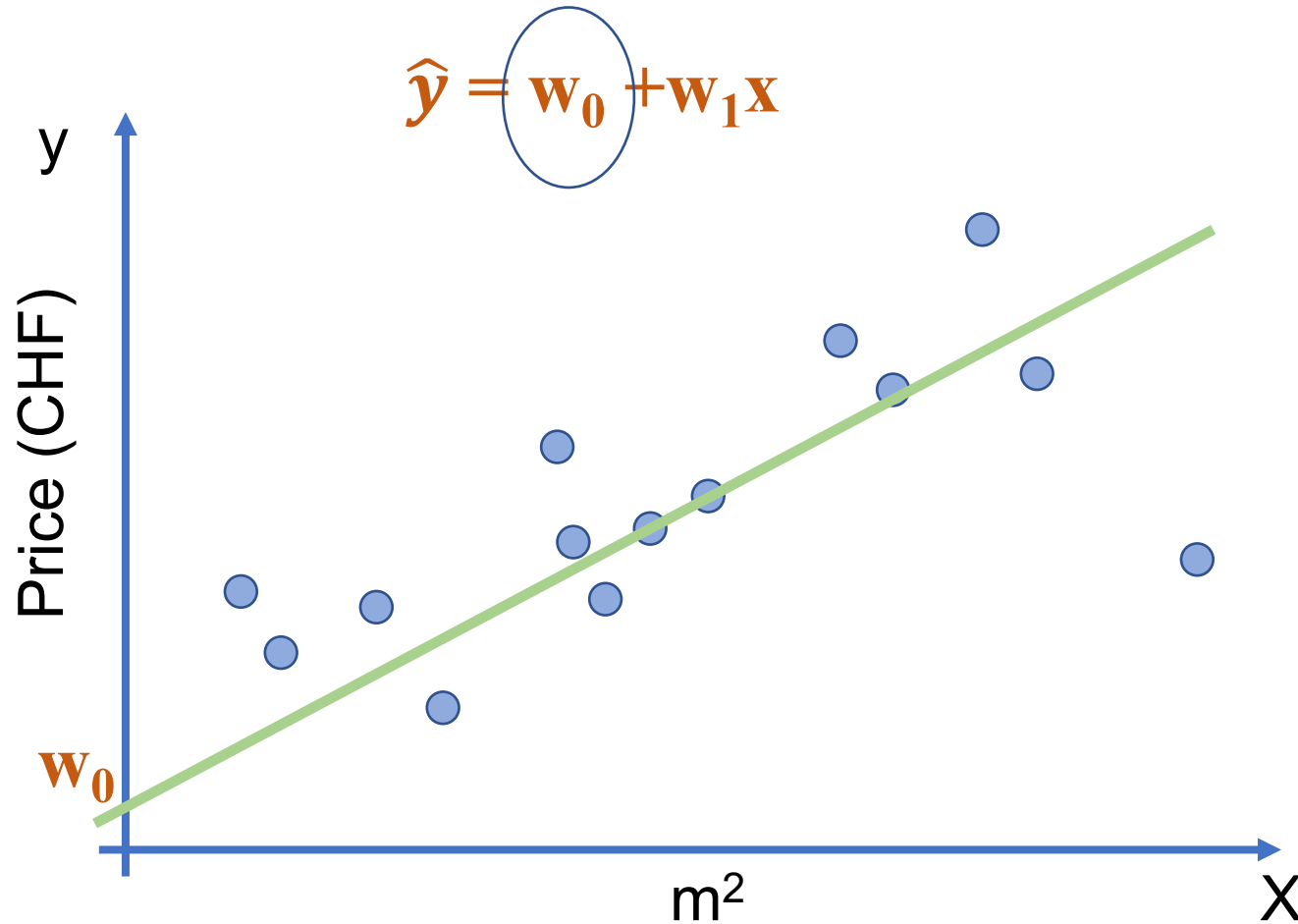
What is the price for a house with area of 100 sq m² ?

$$\begin{aligned}\text{Price} &= -12000 + 8000 * 100 \\ &= 788'000 \text{ CHF}\end{aligned}$$

What is the area for a house sold for 1'200'000?

Interpreting the coefficients

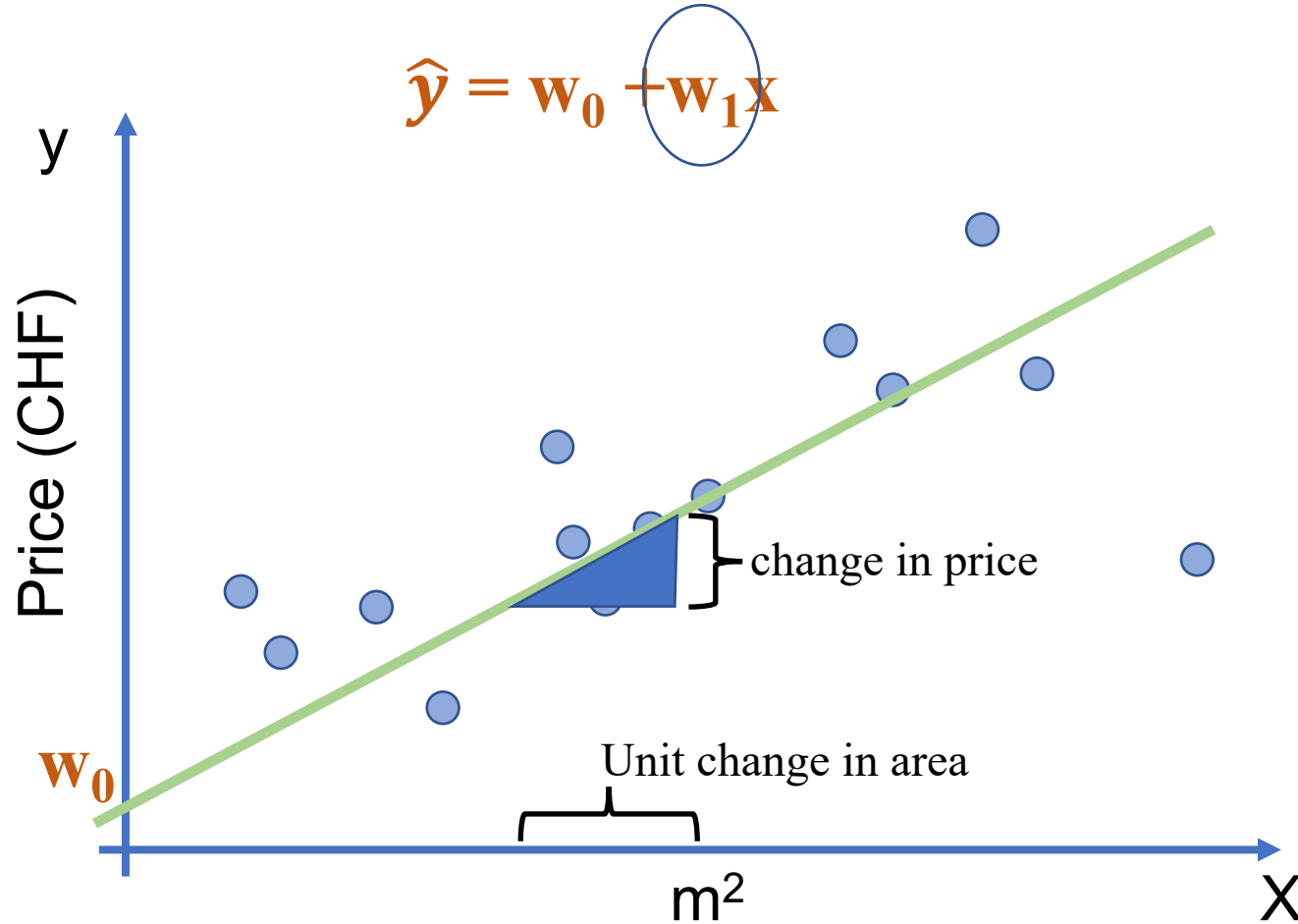
30



Predicted price for house
with zero area
(not very meaningful in
this case)

Interpreting the coefficients

31



Predicted change in house price for 1 unit change in house area.

Exercise:

Compute

$$\hat{y}(101) - \hat{y}(100)$$

Goodness of fit (R^2)

32

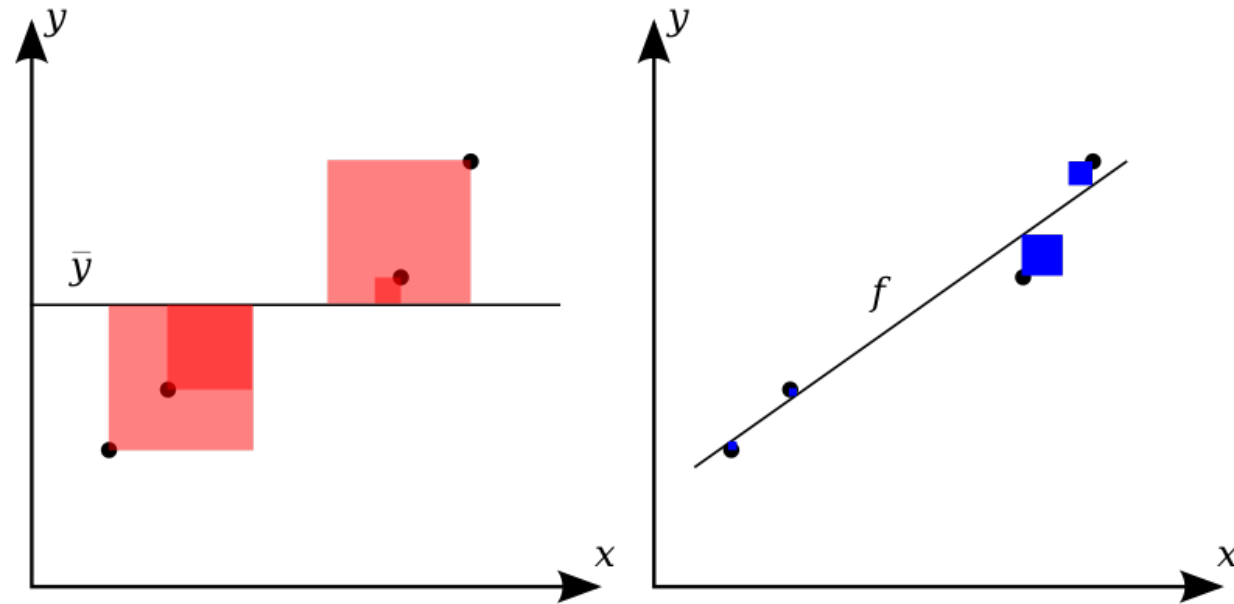
- Typically in statistical analysis we use a number such as R^2 which encodes how much of the data variance is explained by the model.
- $R^2=1$ (perfect model)
- $R^2=0$ (same as baseline, ie avg)

$$R^2 = 1 - \frac{RSS}{TSS}$$

← Explained variance
← Total variance

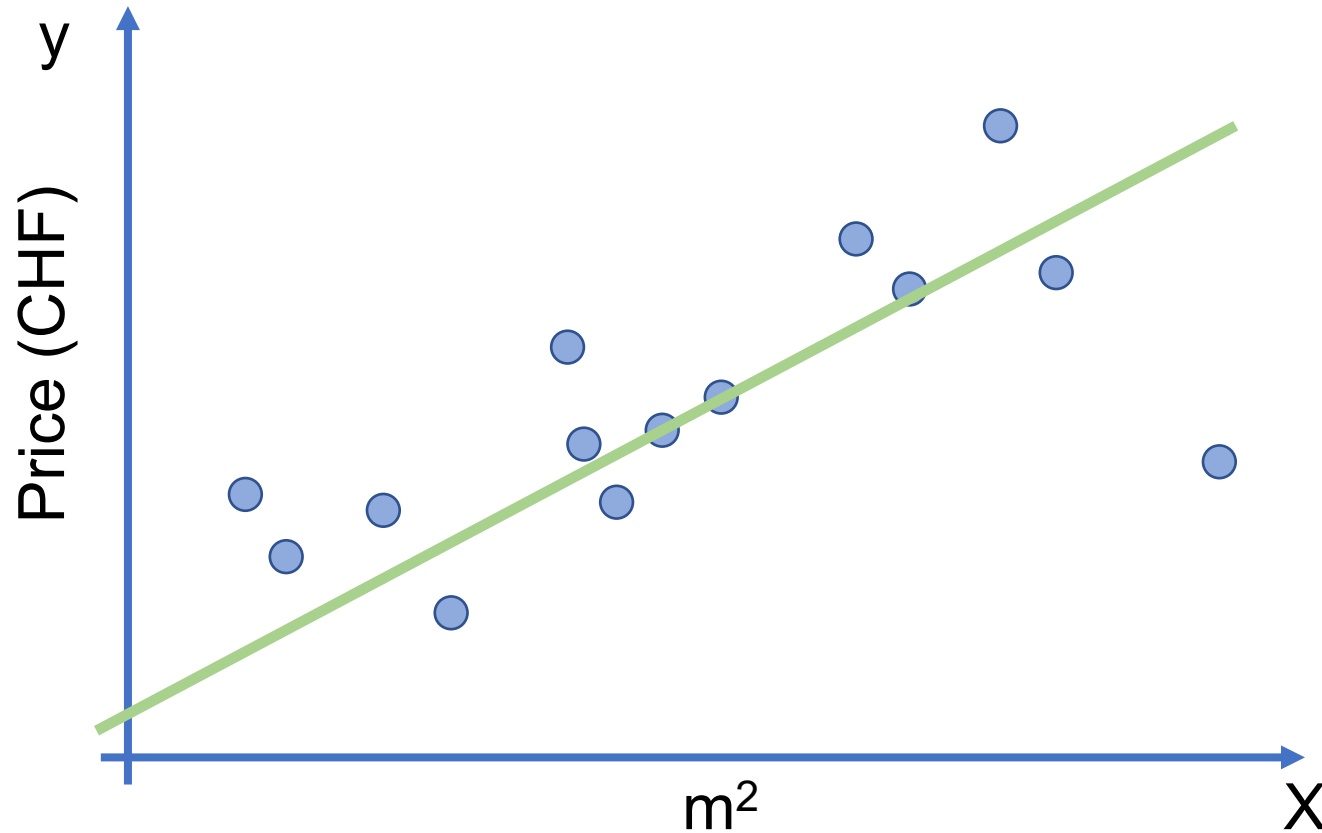
Residual Sums of Squares: $RSS = \sum (y_i - \hat{y}_i)^2$

Total Sums of Squares: $TSS = \sum (y_i - \bar{y})^2$



Goodness of fit (MAE)

33



MAE(Mean Average Error):

$$\frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

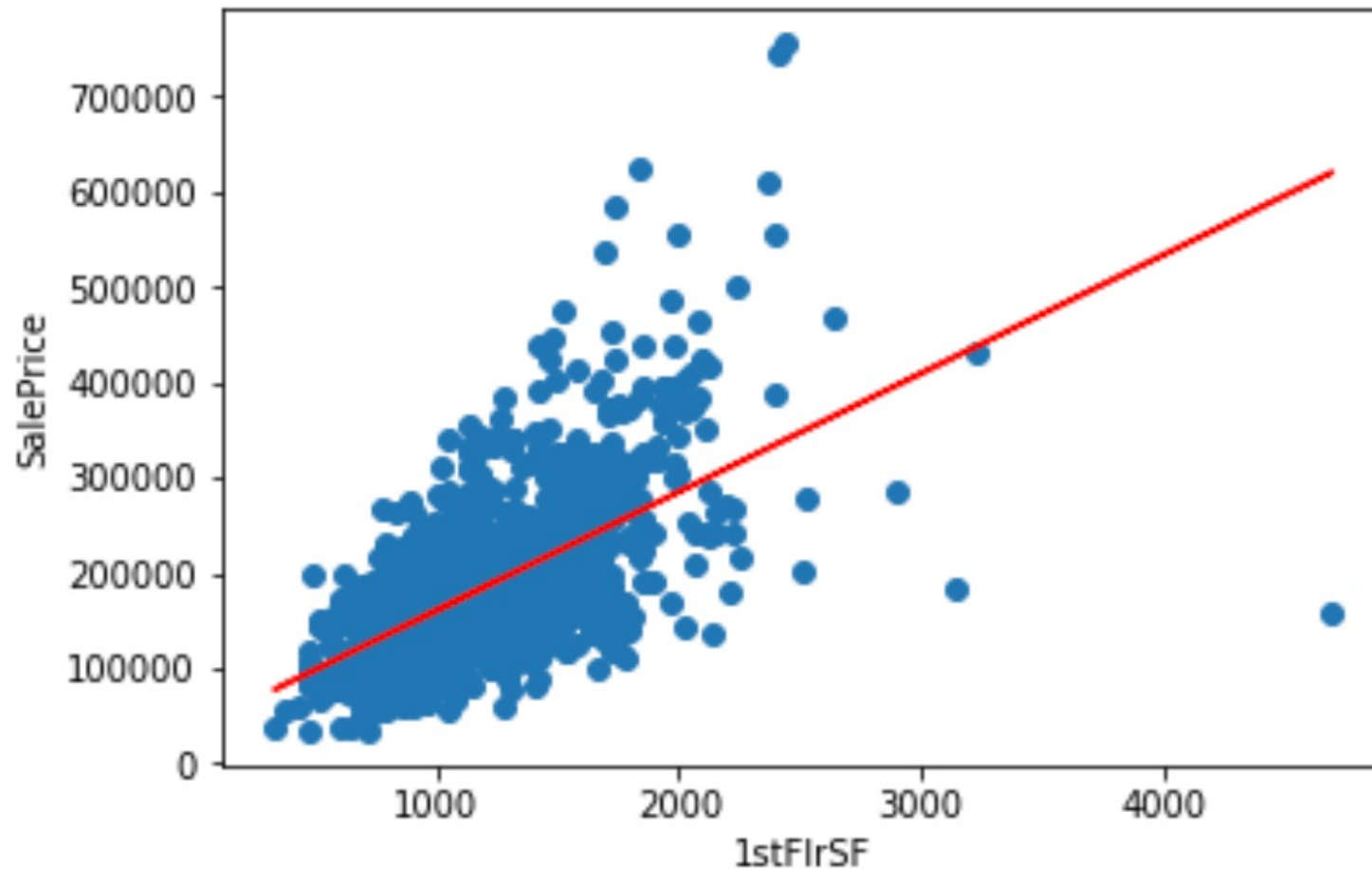
Can provide some more meaningful number (compared to RSS)

Eg if MAE = 25'000 it means that *on average* we are off by 25k CHF in our predictions.

In-class demo – Predicting house prices

34

<https://tinyurl.com/DMML-regression>



In-class Exercise (5 mins)

35

- Do a regression of the features
 - `FullBath` vs `SalePrice`
 - `OverallQual` vs `SalePrice`

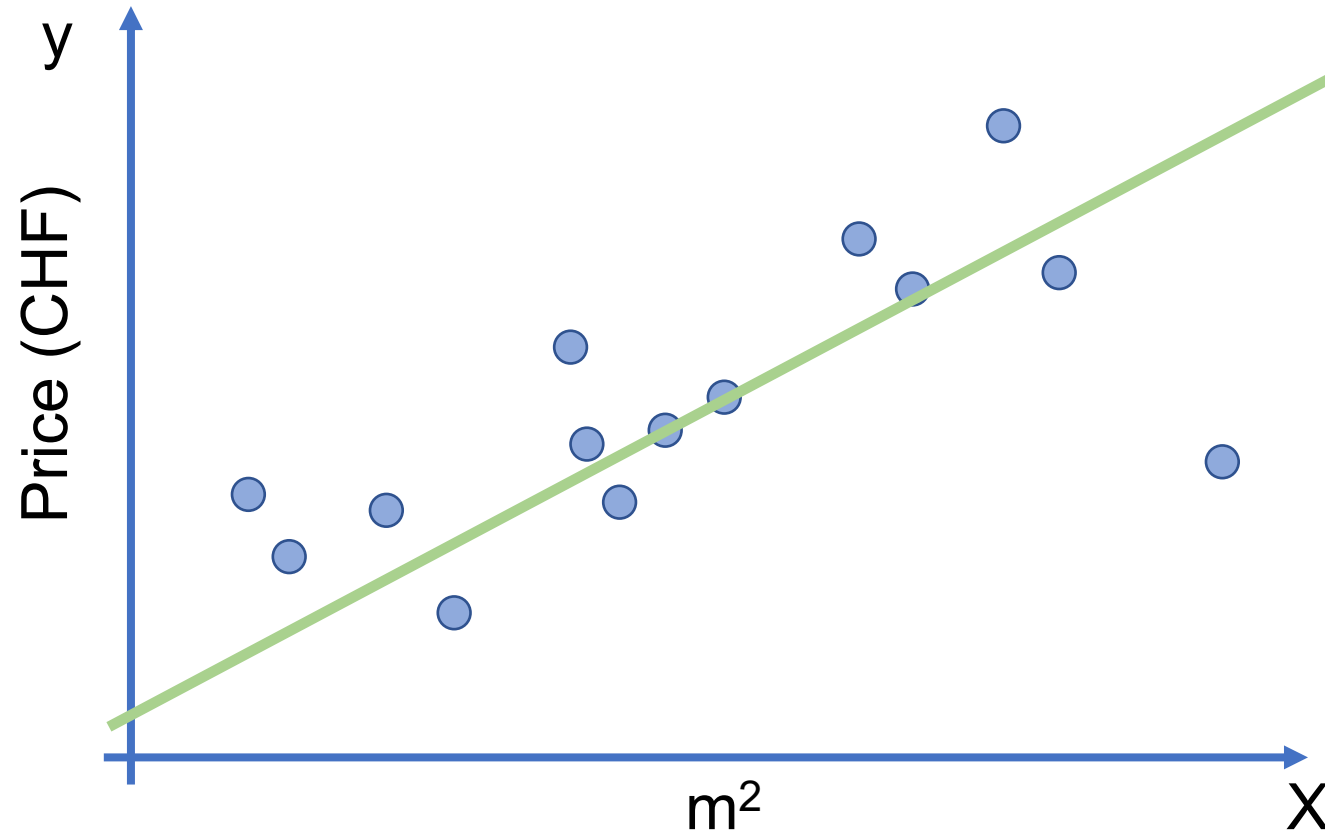
Which has the lowest Mean Average Error (MAE) and R^2 ?

Multiple Regression

Adding more features

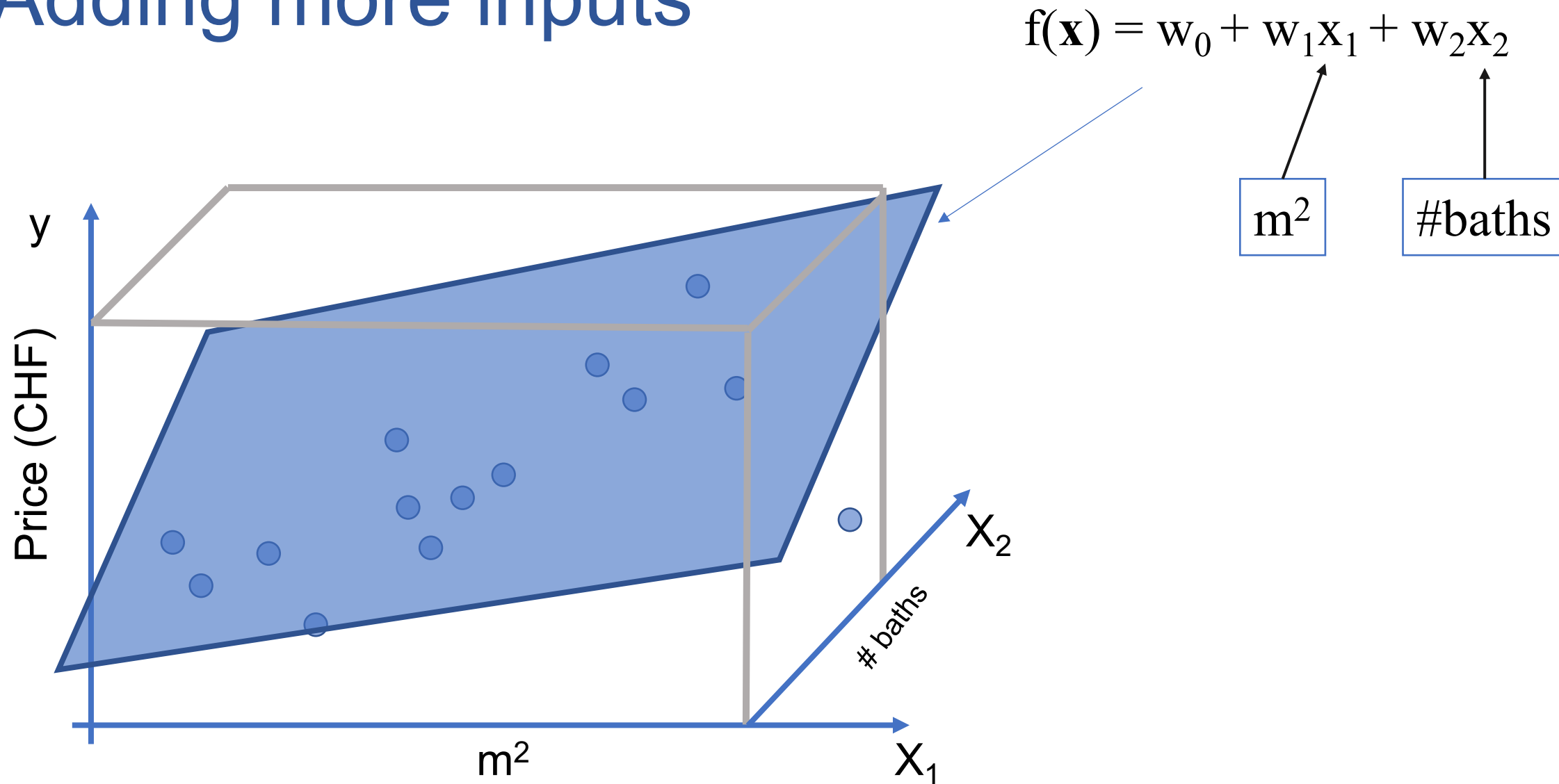
So far we made predictions using only the house size

37



Adding more inputs

38



Many possible inputs

39

- House area
- Lot area
- # bathroom
- # bedrooms
- Year built
- ...

Examples: $n = 4$; dimensions: $d = 4$

40

x_0	House Size x_1	# Rooms x_2	# Bathrooms x_3	Landsize x_4	Price y
1	120	4	2	200	900'000
1	132	5	1	15	1'200'000
1	98	2	1	0	850'000
1	85	2	1	49	970'000

$$X = \begin{bmatrix} 1 & 120 & 4 & 2 & 200 \\ 1 & 132 & 5 & 1 & 15 \\ 1 & 98 & 2 & 1 & 0 \\ 1 & 85 & 2 & 1 & 49 \end{bmatrix}$$

$n \times (d+1)$ array

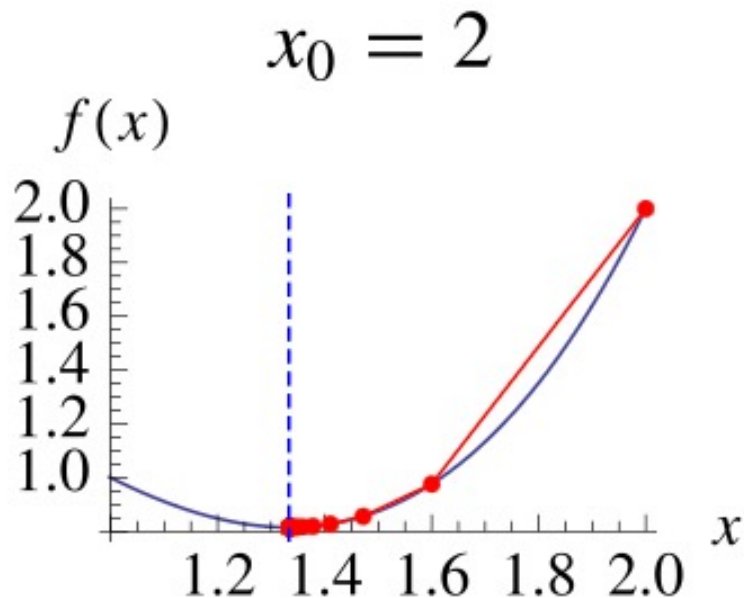
$$y = \begin{bmatrix} 900000 \\ 1200000 \\ 850000 \\ 970000 \end{bmatrix}$$

n -dimensional vector

Steepest (gradient) descent

41

- To find the weights of the linear regression (or of most ML problems) we usually use an approach called gradient descent.



See also here from a demo:

<https://demonstrations.wolfram.com/CurvesOfSteepestDescentFor3DFunctions/>

For those interested in more details, have a look at the advanced (optional) slides.

More regression examples

Can you come up with some more??

Work in Two:

43

- [3mins] Work with your neighbor.
- Find some more interesting regression problems.
 - How would you model it as a regression?
 - What are the features that influence the target variable?
- [3mins] We discuss.

Determine salary

44

- How much salary y will you get after you graduate?
- Depends on X=
 - total grade,
 - #internships,
 - years of past experience,
 - #friends on Facebook,
 - #people you connect with in LinkedIn,
 - ...

$$\hat{y} = w_0 + w_1 \text{grade} + w_2 \text{experience} + w_3 \text{Facebook_friends} + \dots$$

weight

but cannot compare the weight unless i normalize the data.

Stock price prediction

45

Depends on:

- Last day's *value*
- Last week's price *slope*
- How many events in the news about the company
- # times Donald Trump mentioned the company
- ...

Facebook, Inc. (FB)

NasdaqGS - NasdaqGS Real Time Price. Currency in USD

★ In watchlist

185.67 +0.10 (+0.05%) **183.70** -1.97 (-1.06%)

At close: August 30 4:00PM EDT

Pre-Market: 5:42AM EDT

Buy

Sell

Summary

Chart

Conversations

Statistics

Historical Data

Profile

Financials

Analysis

Options

Holders

Sustainability

Previous Close	185.57	Market Cap	518.833B
Open	186.78	Beta (3Y Monthly)	1.25
Bid	183.73 x 1400	PE Ratio (TTM)	31.40
Ask	184.42 x 1300	EPS (TTM)	5.91
Day's Range	183.46 - 186.80	Earnings Date	Oct 28, 2019 - Nov 1, 2019
52 Week Range	123.02 - 208.66	Forward Dividend & Yield	N/A (N/A)
Volume	10,785,722	Ex-Dividend Date	N/A
Avg. Volume	16,752,906	1y Target Est	232.33

Trade prices are not sourced from all markets



Analyst Recommendation
by Argus Research

BUY

Fair Value
[View details](#)

Undervalued



All

Short Term

Mid Term

Long Term

[View more ideas](#)

Predict number of retweets

- How many people will retweet your post?
- Depends on:
 - # followers
 - # hashtags in post
 - Popularity of hashtags
 - Text in post
 - # images in post
 - ...



The image shows a screenshot of Donald J. Trump's Twitter profile. The profile header includes a circular profile picture of Trump, a background image of a crowd wearing red hats, and statistics: 44K Tweets, 47 Following, 63.9M Followers, 7 Likes, and 6 Moments. The bio identifies him as the 45th President of the United States, located in Washington, DC, with links to his Instagram and a note that he joined in March 2009. Below the bio is a gallery of 3,315 photos and videos. The main content area shows a tweet from the National Hurricane Center (@NHC_Atlantic) retweeted by Donald J. Trump. The tweet, posted 7 hours ago, contains key messages for Hurricane Dorian, including warnings about storm surge, life-threatening winds, and heavy rains. It includes a map of the hurricane's path and a link to hurricanes.gov for more information. At the bottom, there are engagement metrics: 598 replies, 2.4K retweets, and 6.4K likes. A partial tweet from Donald J. Trump is visible at the very bottom, mentioning a trade agreement.

Donald J. Trump 
@realDonaldTrump

45th President of the United States of America 

Washington, DC

[Instagram.com/realDonaldTrump](https://www.instagram.com/realDonaldTrump)

Joined March 2009

[3,315 Photos and videos](#)

Donald J. Trump Retweeted

National Hurricane Center  @NHC_Atlantic · 7h

Here are the 11 PM EDT Monday, September 2 Key Messages [#Dorian](#). For more information, visit hurricanes.gov.

Key Messages for Hurricane Dorian
Advisory 39: 11:00 PM EDT Mon Sep 02, 2019

1. Devastating winds and storm surge will continue to affect Grand Bahama Island for several more hours. Everyone there should remain in shelter.
2. Life-threatening storm surge and dangerous hurricane-force winds are expected along portions of the Florida east coast and the coasts of Georgia and South Carolina, regardless of the exact track of Dorian's center. Water levels could begin to rise well in advance of the arrival of strong winds. Residents in these areas should follow advice given by local emergency officials.
3. The risk of life-threatening storm surge and hurricane-force winds continues to increase along the coast of North Carolina. Residents in these areas should follow advice given by local emergency officials.
4. Heavy rains, capable of producing life-threatening flash floods, are expected over northern portions of the Bahamas and coastal sections of the southeast and lower mid-Atlantic regions of the United States through Friday.

For more information go to hurricanes.gov

598 2.4K 6.4K

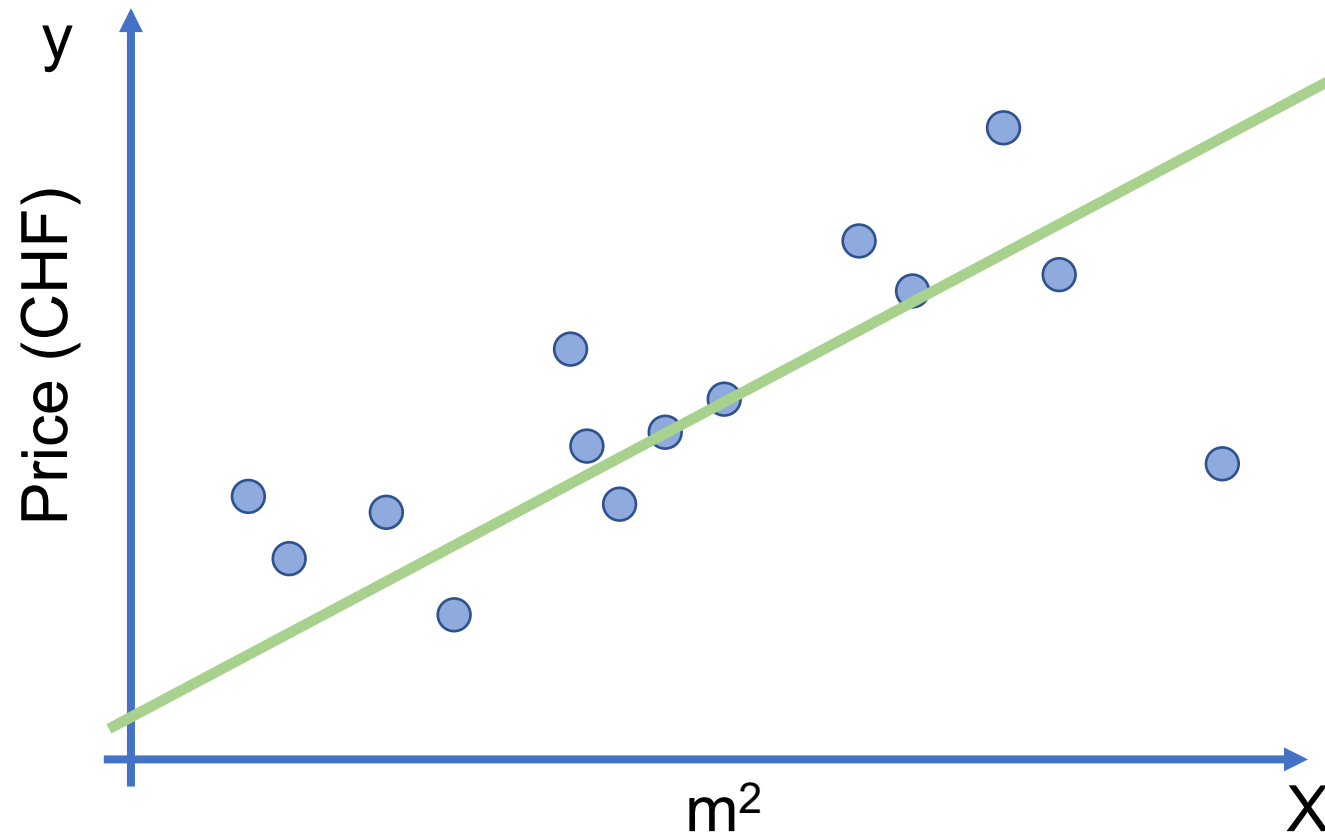
Donald J. Trump  @realDonaldTrump · 7h

...Trade Agreement." [@business](#) [@ChuckGrassley](#) [@joniern](#)
[@BenSasse](#) Making great progress for our Farmers. Appro

Building more complex models

Our linear model

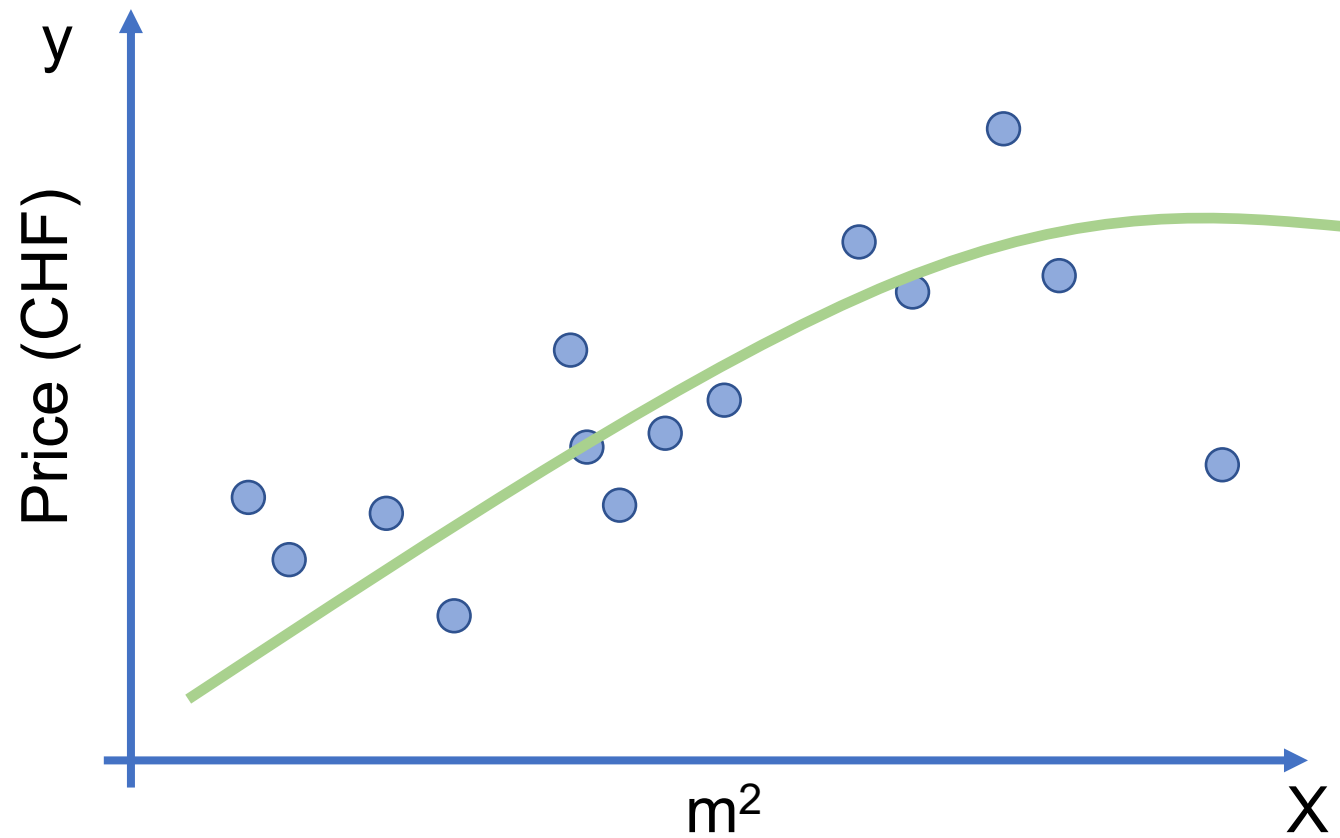
48



$$\text{RSS} = 130$$

A quadratic model

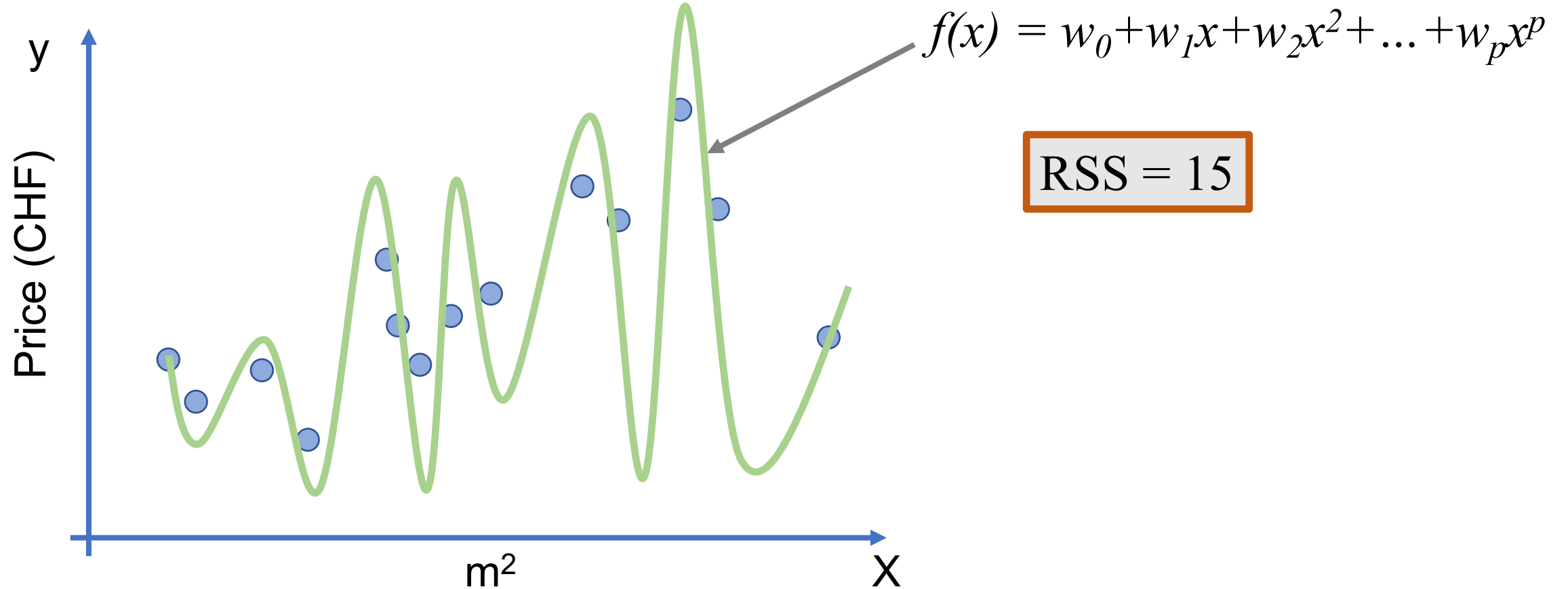
49



$$RSS = 110$$

A higher-order polynomial

50



Example: Polynomial regression

51

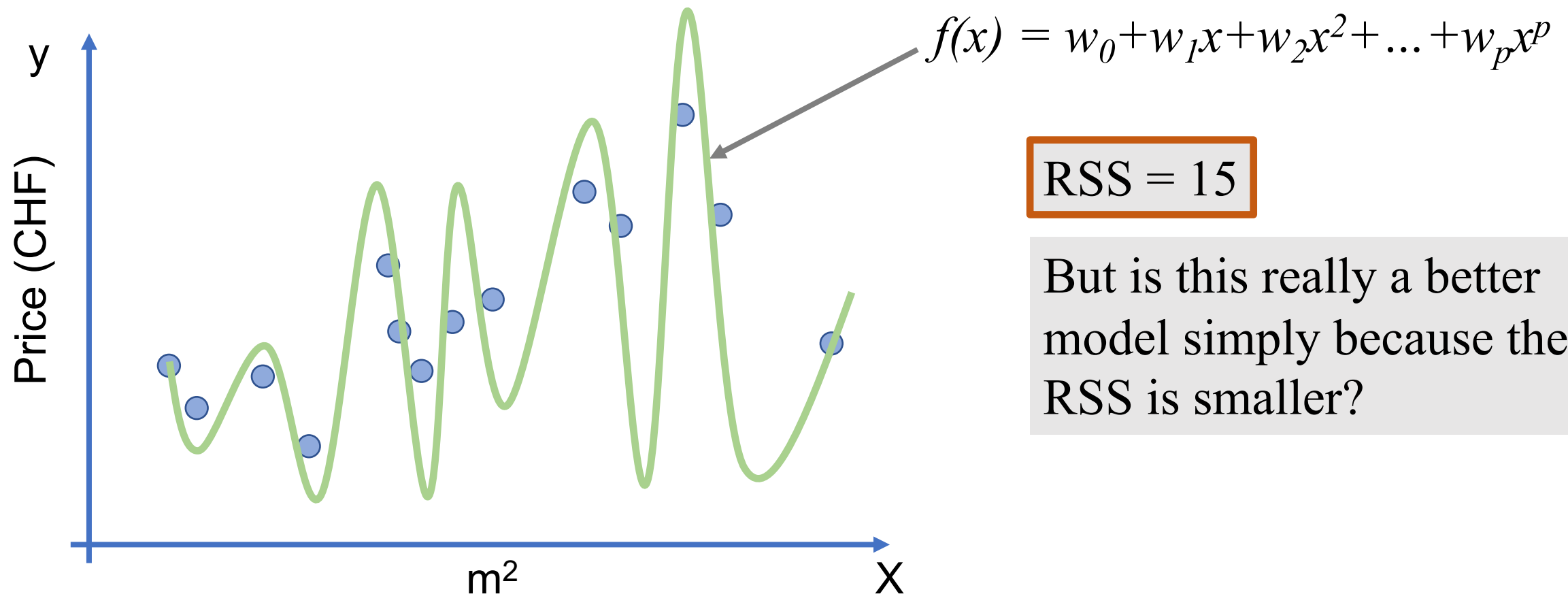
x_0	House Size x	House Size ² x^2	Price y
1	120	120^2	900'000
1	132	132^2	1'200'000
1	98	98^2	850'000
1	85	85^2	970'000

$$X = \begin{bmatrix} 1 & 120 & 120^2 \\ 1 & 132 & 132^2 \\ 1 & 98 & 98^2 \\ 1 & 85 & 98^2 \end{bmatrix}$$

$$y = \begin{bmatrix} 900000 \\ 1200000 \\ 850000 \\ 970000 \end{bmatrix}$$

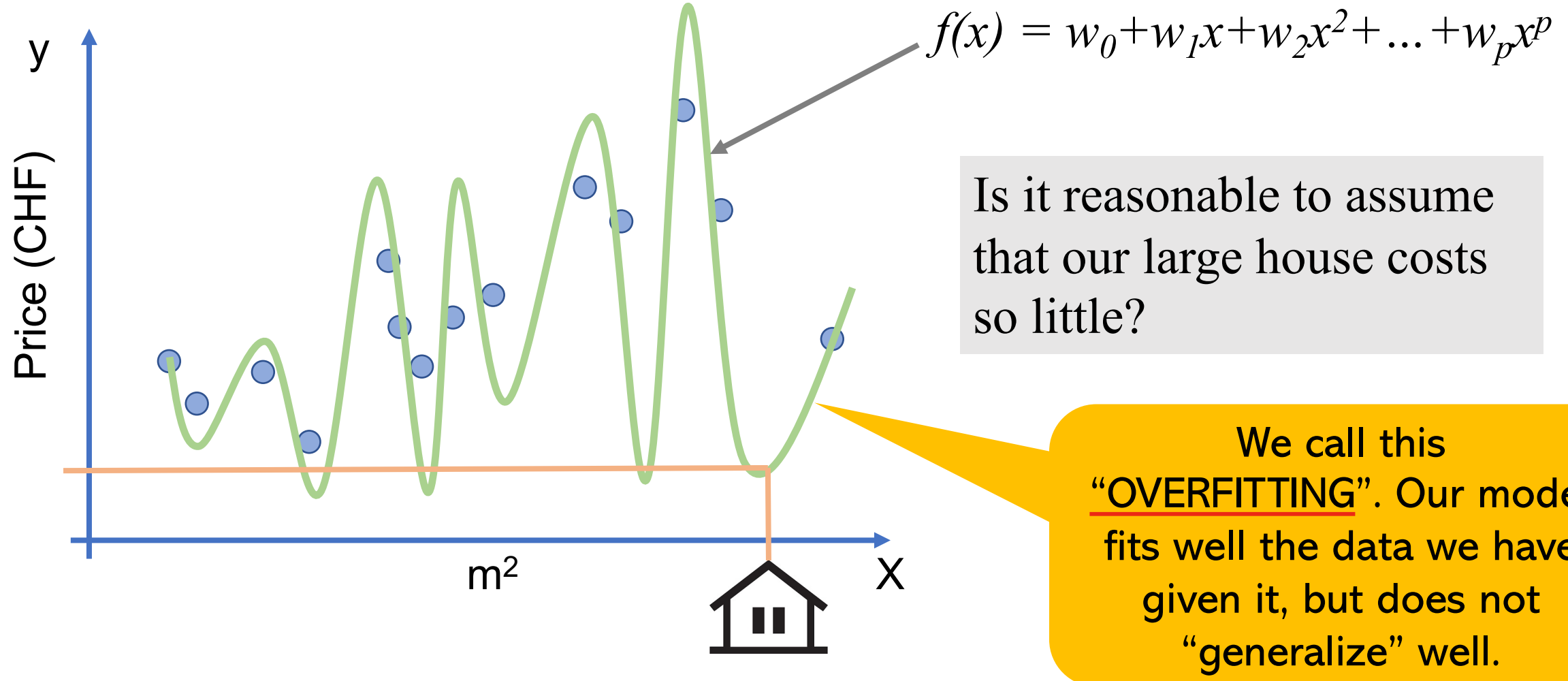
A higher-order polynomial

52



A higher-order polynomial

53



Overfitting

54

“*Overfitting* is the tendency of data mining procedures to tailor models to the training data, at the expense of generalization to previously unseen data points”

See also [Chapter 5](#) of “Data Science for Business” book.

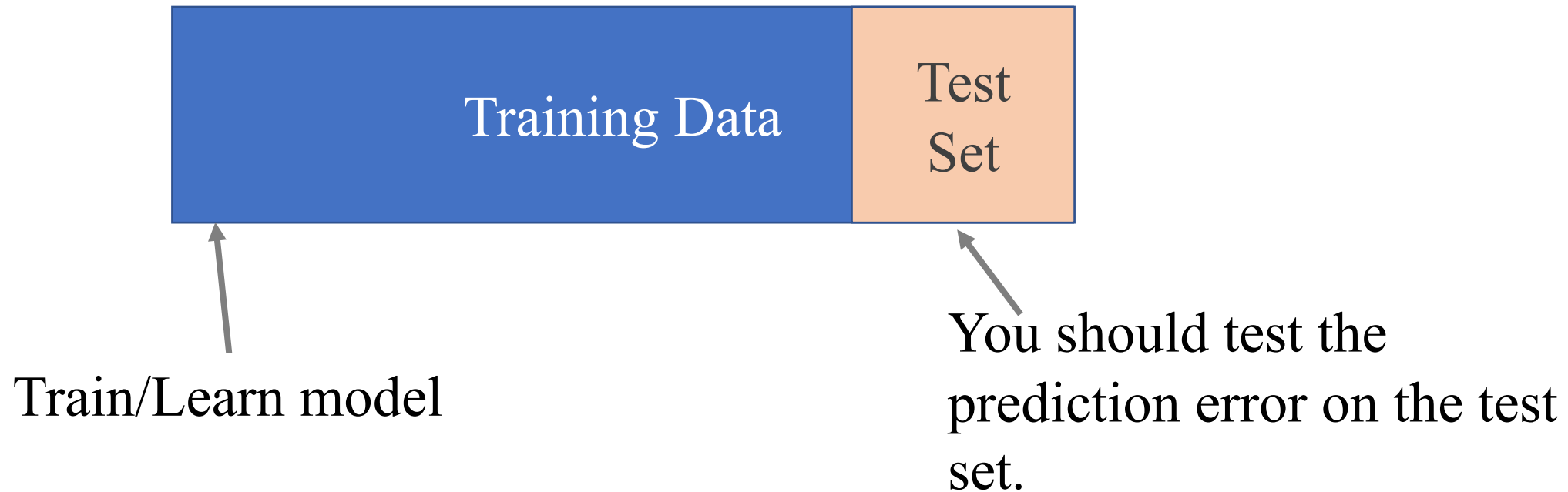
Evaluation

How do we evaluate how good our model is?

Evaluation of your model

56

- This is the **most important slide** in this class, and where most people make the **most errors** when modeling data!



Evaluate error on data you haven't seen!

57

- We would like to have good predictions, but we don't know what the future brings (nobody does!).

Simulate predictions:

- Remove observations (create training/test set)
- Fit model on remaining (use **training set**)
- Predict on held-out observations (use **test set**)

Whole dataset

House area	Lot area	# baths	Sale Price
120	500	...	930k
65	350	...	705k
154	0	...	2010k
220	0	...	3000k
65	15	...	350k
85	35	...	810k
122	0	...	1200k

Evaluate error on data you haven't seen!

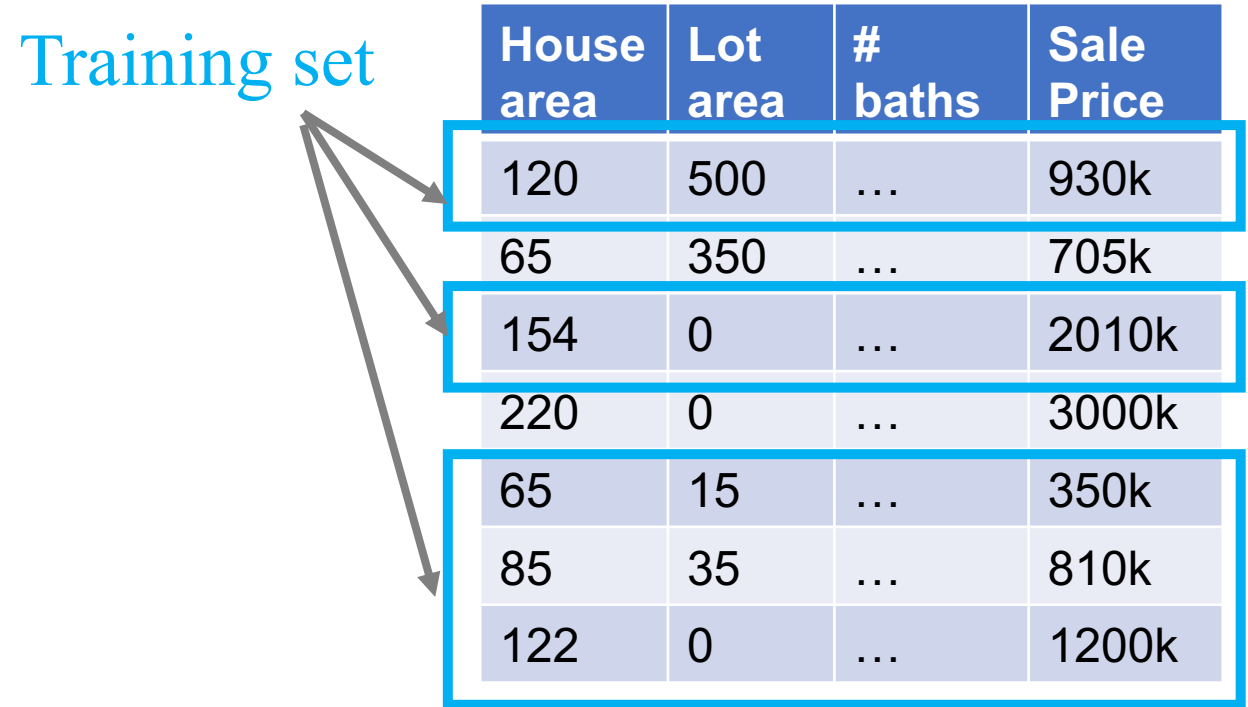
58

- We would like to have good predictions, but we don't know what the future brings (nobody does!).

Simulate predictions:

- Remove observations (create training/test set)
- Fit model on remaining (use **training set**)
- Predict on held-out observations (use **test set**)

Training set



House area	Lot area	# baths	Sale Price
120	500	...	930k
65	350	...	705k
154	0	...	2010k
220	0	...	3000k
65	15	...	350k
85	35	...	810k
122	0	...	1200k

Evaluate error on data you haven't seen!

59

- We would like to have good predictions, but we don't know what the future brings (nobody does!).

Simulate predictions:

- Remove observations (create training/test set)
- Fit model on remaining (use **training set**)
- Predict on held-out observations (use **test set**)

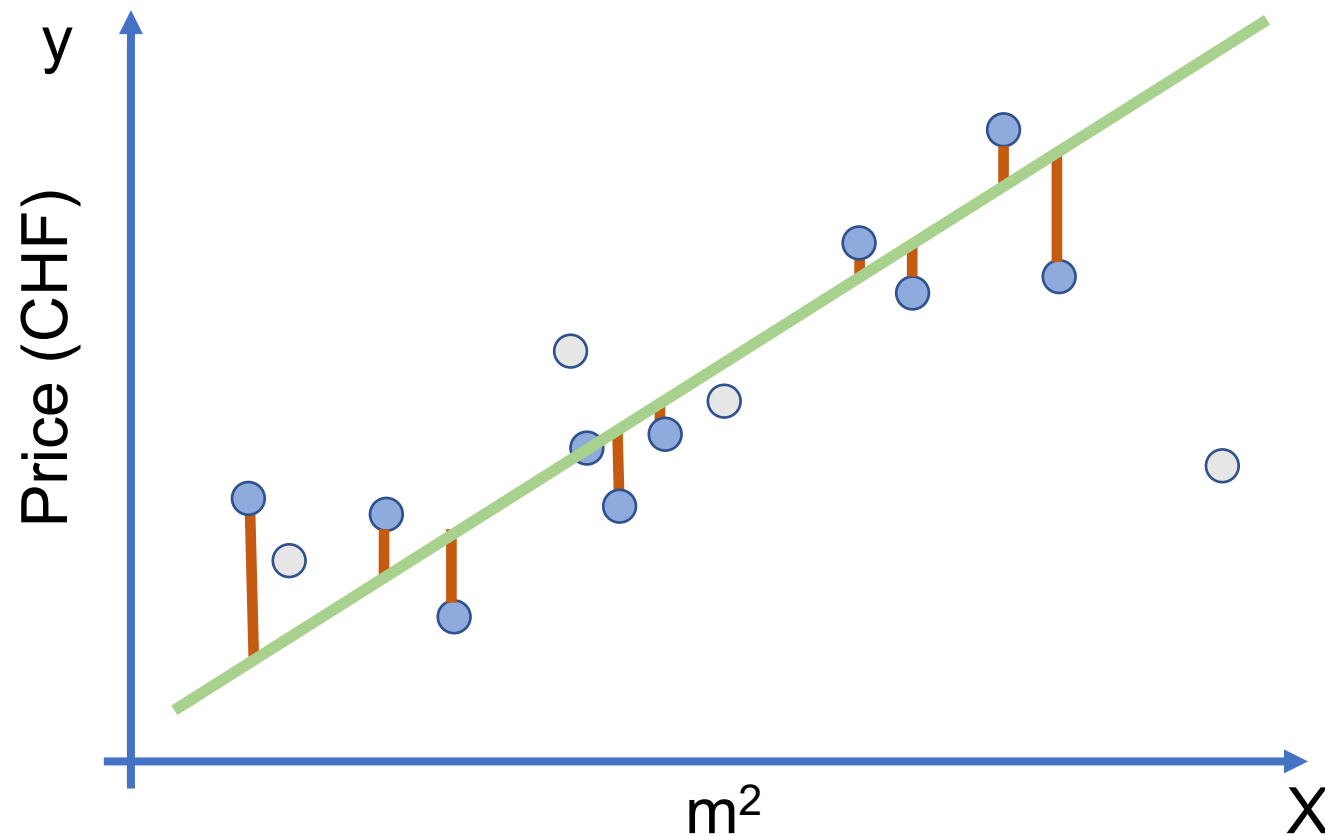
test set

House area	Lot area	# baths	Sale Price
120	500	...	930k
65	350	...	705k
154	0	...	2010k
220	0	...	3000k
65	15	...	350k
85	35	...	810k
122	0	...	1200k

Training Error

60

Use only training set

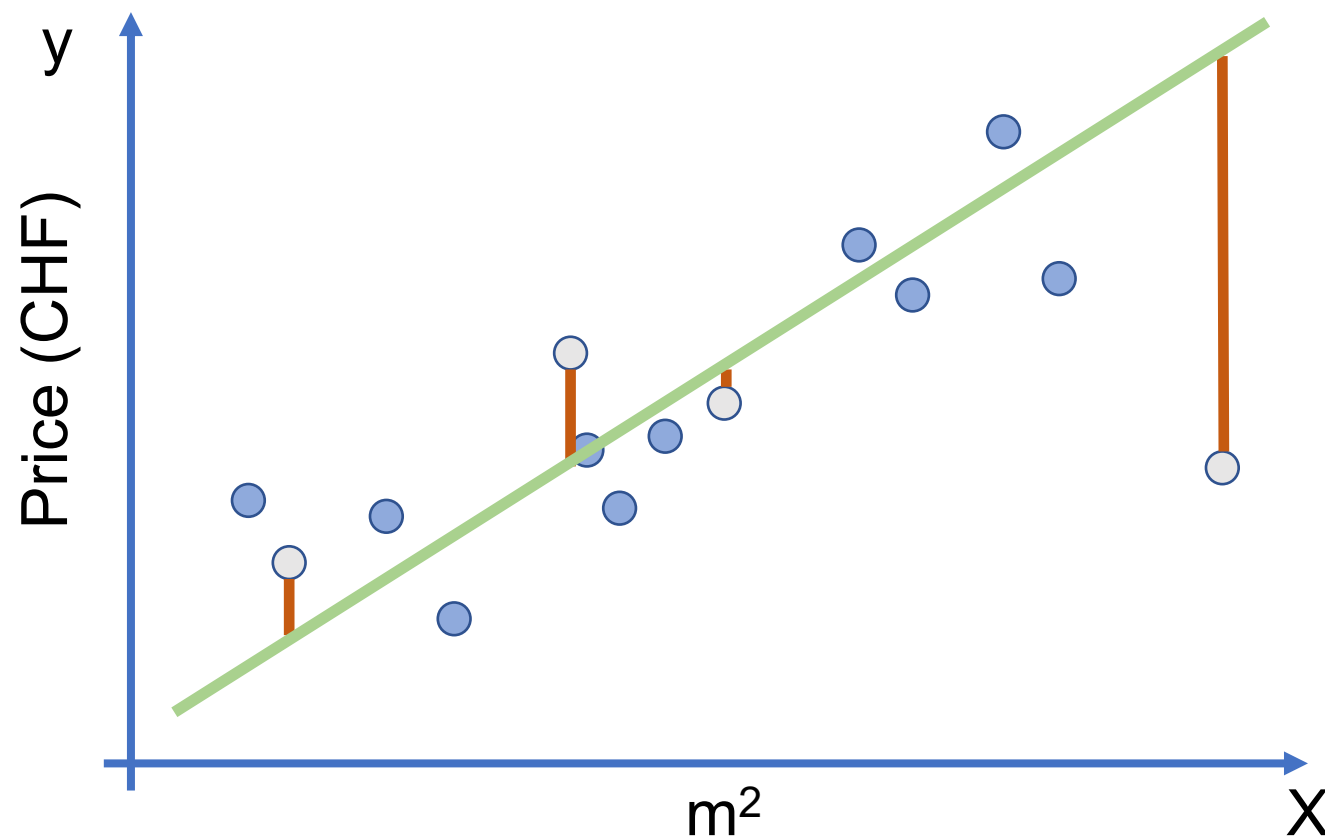


$$\begin{aligned} \text{Training error}(w) = & (y_{\text{training}1} - f_w(x_{\text{training}1}))^2 \\ & + (y_{\text{training}2} - f_w(x_{\text{training}2}))^2 \\ & + \dots \\ & [\text{for all training examples}] \end{aligned}$$

Test Error (a more realistic estimate of the prediction error)

61

Use only test set



$$\begin{aligned} \text{Test error}(w) = & (y_{\text{test1}} - f_w(x_{\text{test1}}))^2 \\ & + (y_{\text{test2}} - f_w(x_{\text{test2}}))^2 \\ & + \dots \\ & [\text{for all test examples}] \end{aligned}$$

*Assess prediction error
using only the test set*

80-20 split

62

- Typically we use an 80-20 split for train/test sets.
 - This means we use
 - 80% of the dataset for training the model (learning the parameters w)
 - 20% as the test set for predicting the model accuracy (or error)
- A. First we **shuffle** the rows of the dataset.
- B. We select the first 80% of the rows → train set
- C. We select the last 20% of the rows → test set



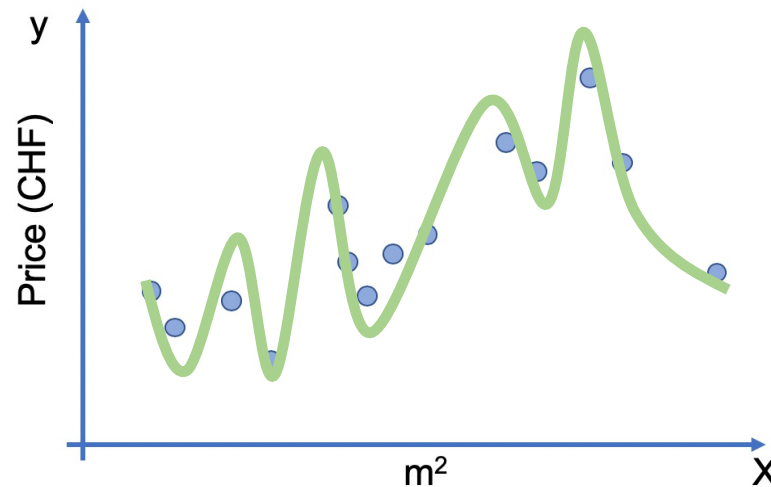
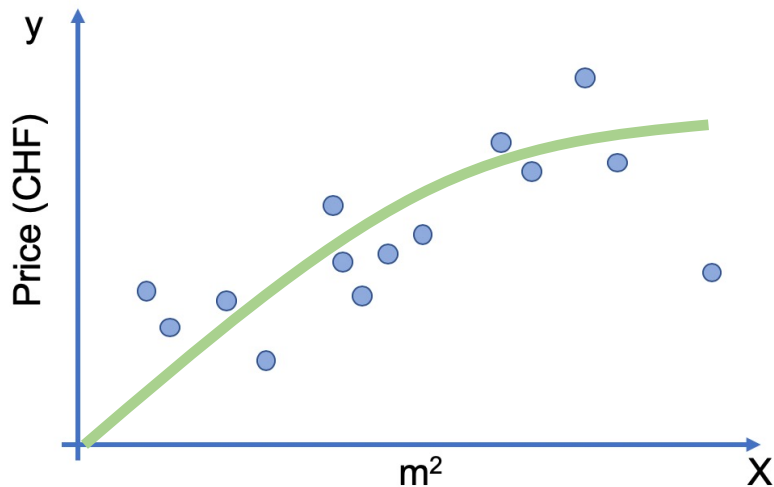
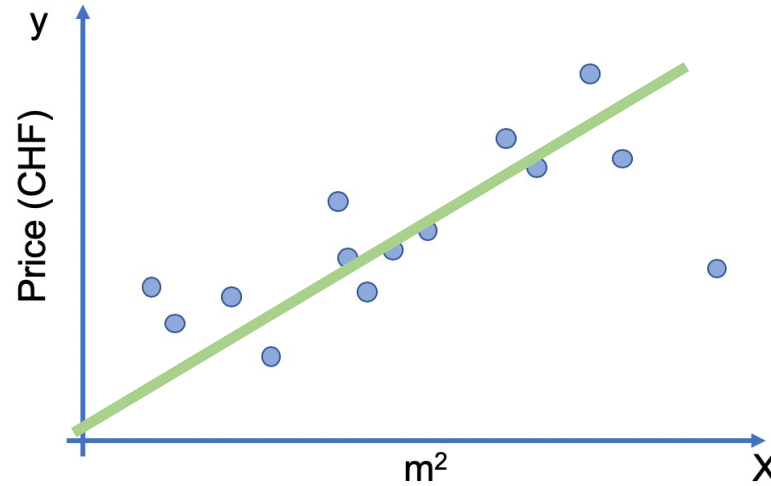
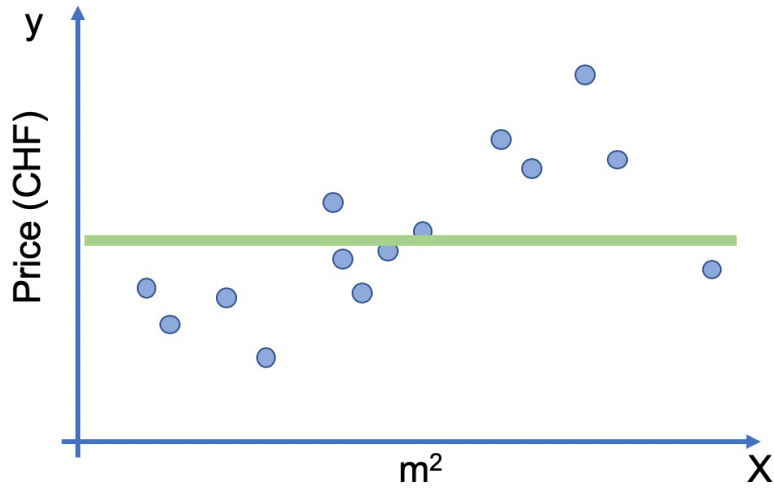
How to evaluate which model is best?

63

- Imagine we are given many models. How do we find which one we should choose?
- We use again the train/test split and found out which behaves the best for the **test data**.

How to evaluate which model is best?

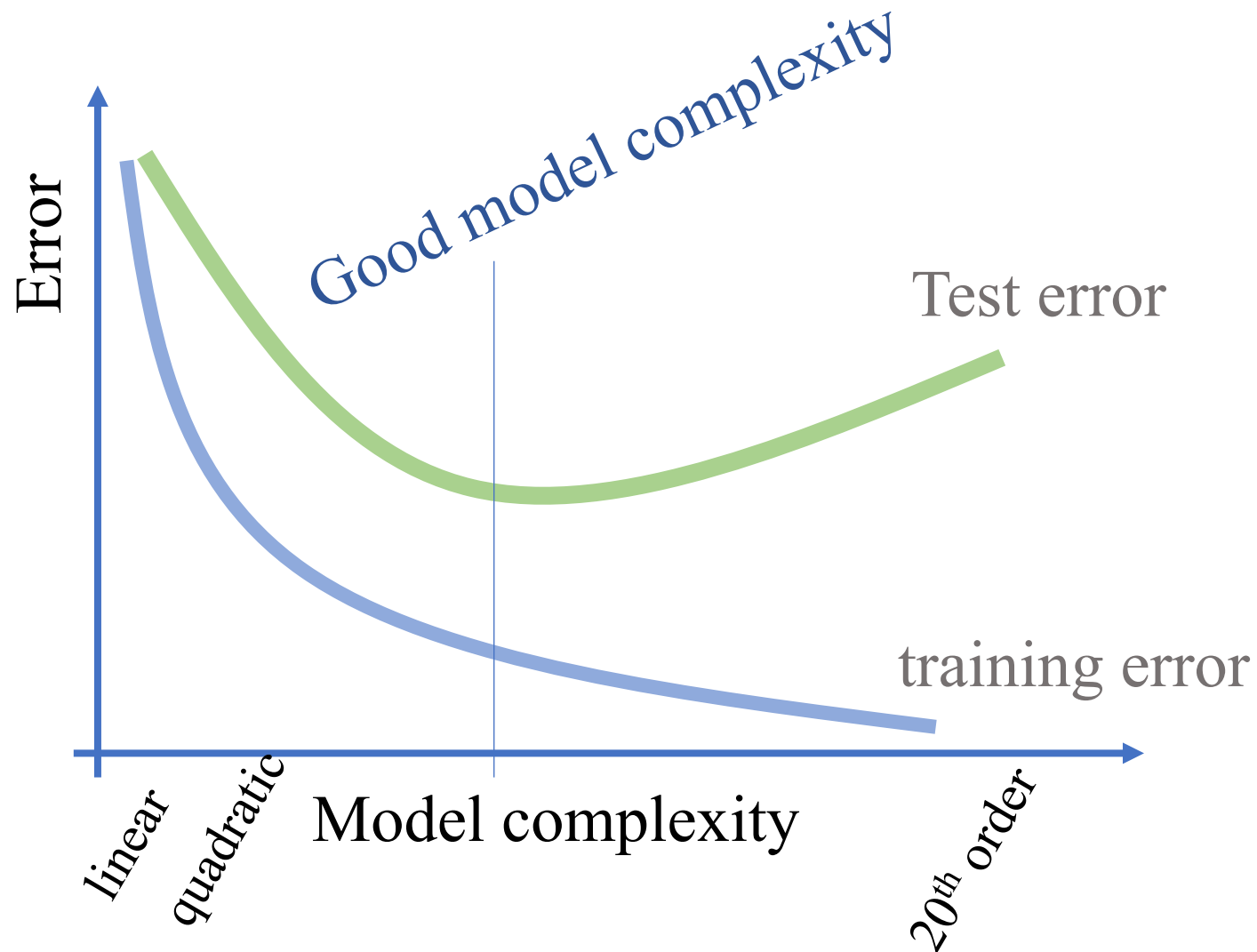
64



To select which of those models is the “best”, means which has the lowest error in the **test data**, because this means that it generalizes well.

Train/Test Curves

65



We pick the model that has the lowest test error. The training error will (almost) always reduce with more complex models. **You want the simplest possible model, with the least test error.**

In-class exercise: Income vs Height

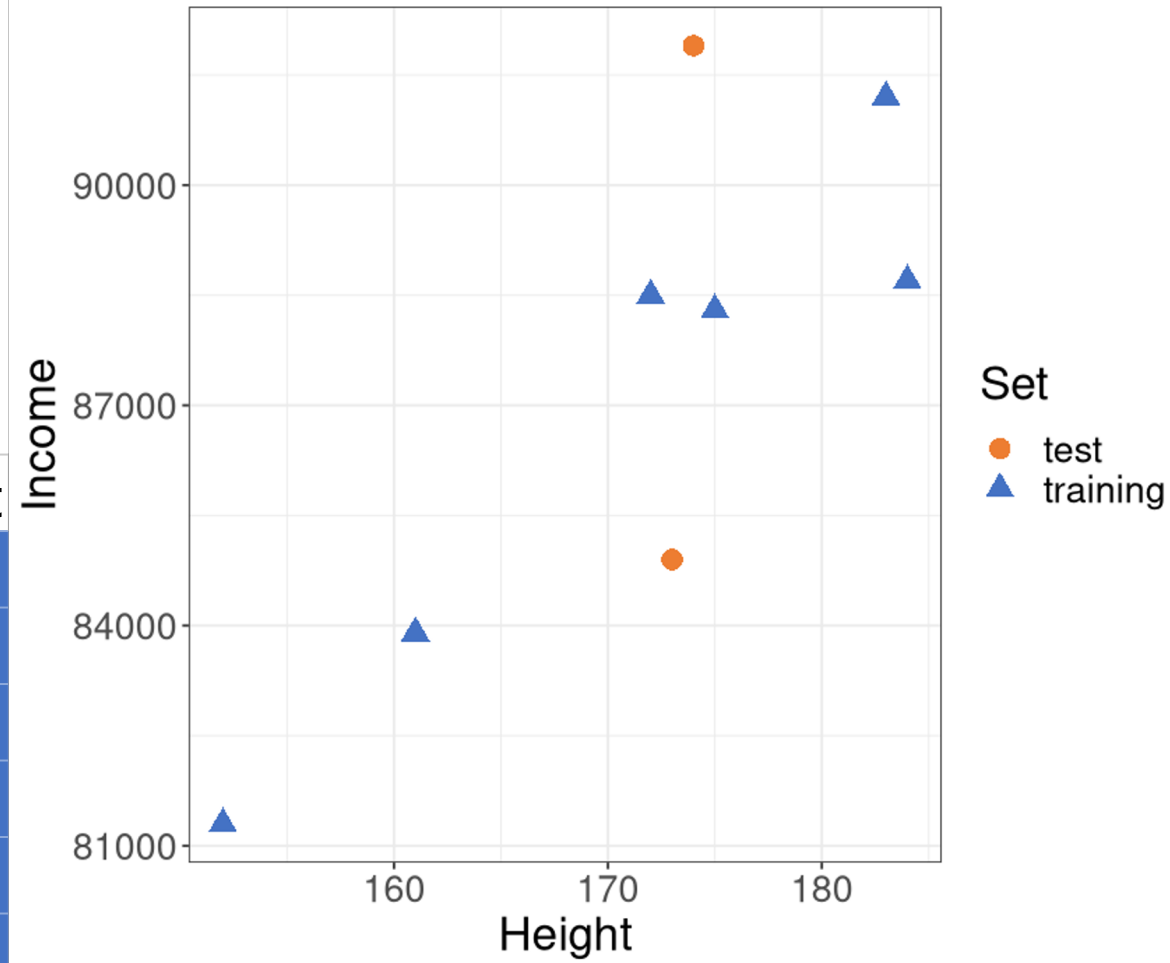
66

We have the following data on income vs height of a person.

1. Compute the MAE of the train and the test data on this model.

$$f(x) = 20'000 + 400 x$$

Income	Height
91200	183
83900	161
88500	172
88300	175
81300	152
88700	184
84900	173
91900	174



In-class exercise: Income vs Height

67

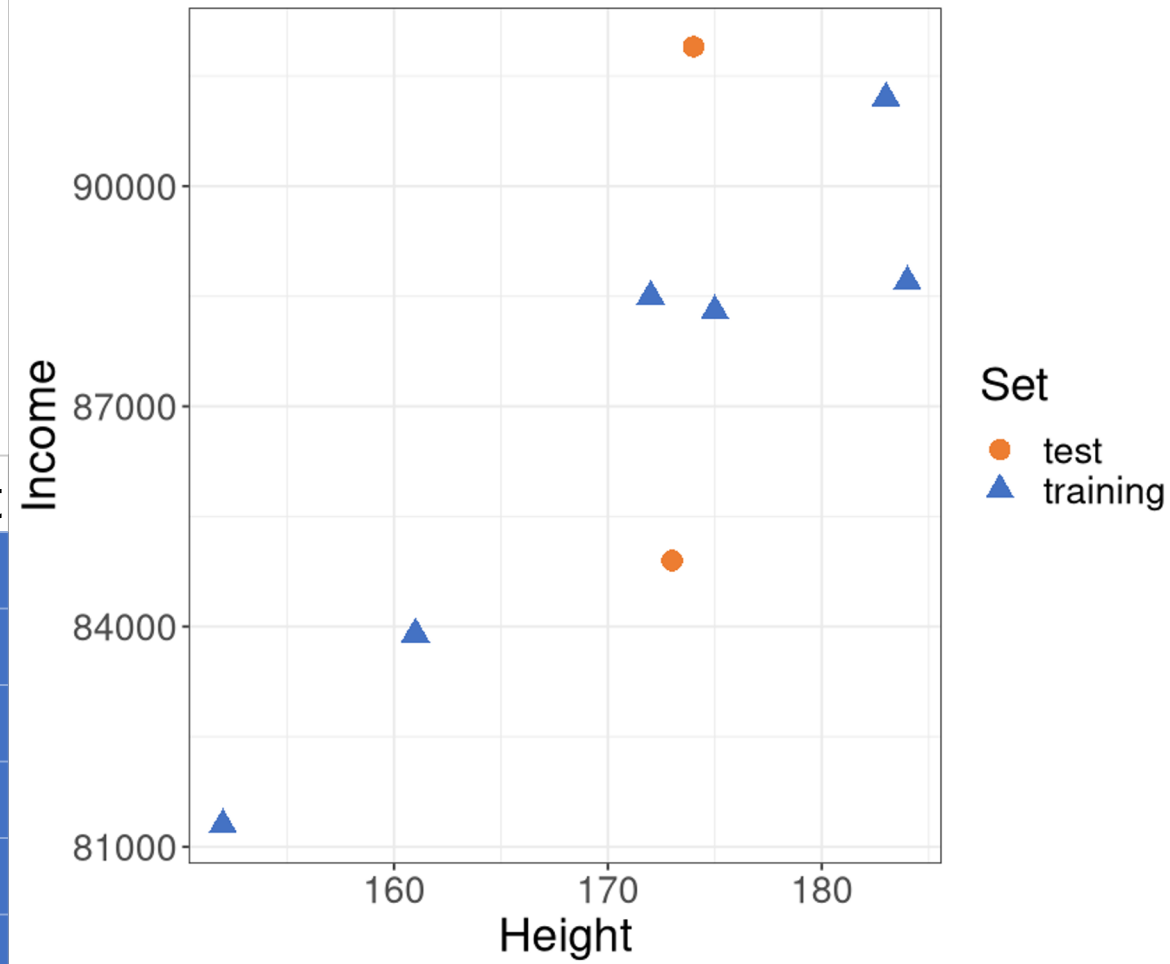
We have the following data on income vs height of a person.

1. Compute the MAE of the train and the test data on this model.

$$f(x) = 20'000 + 400 x$$

2. Given this new model, $f(x) = 38'600 + 290 x$, which one would you pick and why?

Income	Height
91200	183
83900	161
88500	172
88300	175
81300	152
88700	184
84900	173
91900	174



Cross-Validation

More accurate estimates of the prediction error

Creating many test sets

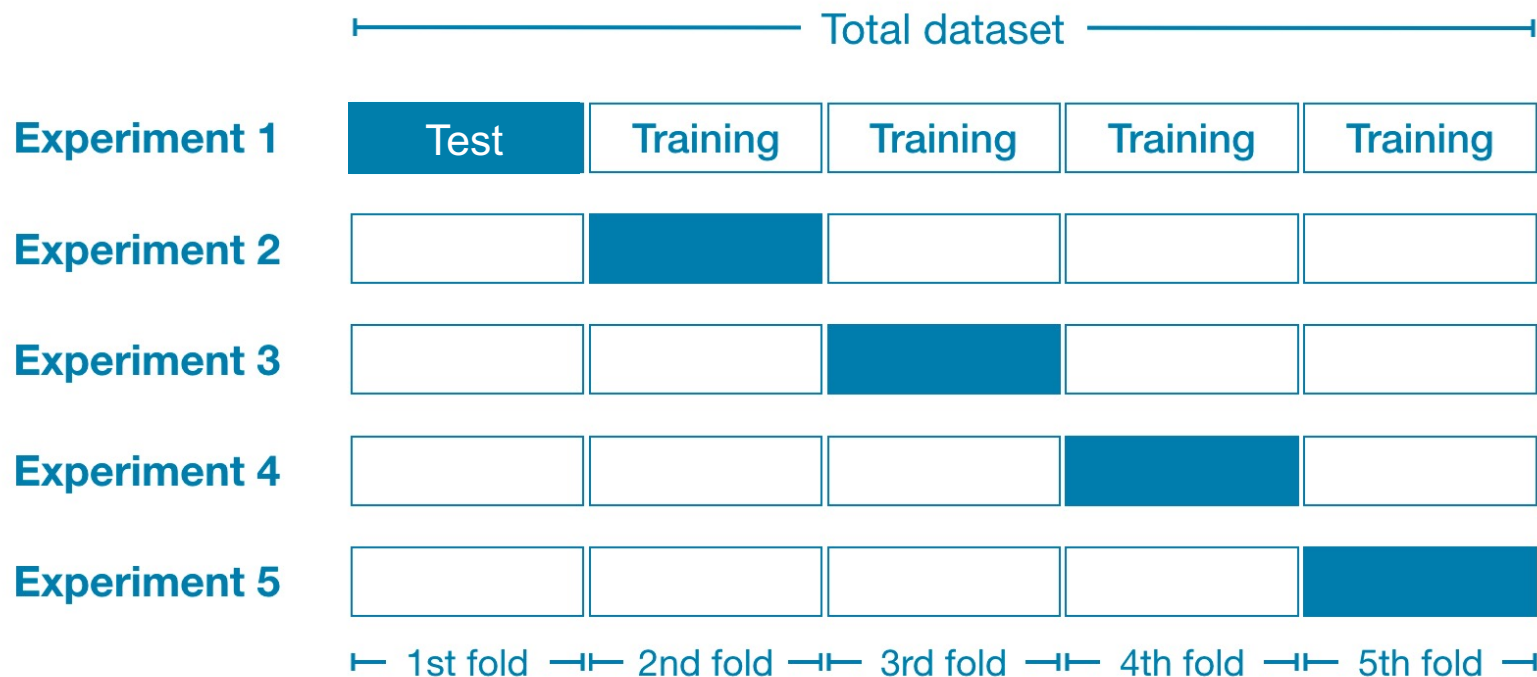
69

- If we use only **one test set**, then we may get *good or bad* prediction results depending on what the test set was.
- We can get a more realistic estimate of the prediction error by using **many test sets**, which will yield a better measure of model quality.

K-fold cross validation (CV)

70

- Typically we use $K=5$ or 10 .
- $\text{Test_Error} = \frac{1}{K} \sum_1^K \text{test_error}_{\text{fold}-K}$
- If $K=\text{size of dataset}$, then we have a “leave-one-out” CV. It provides the most accurate estimate of model quality but it is also the most costly.



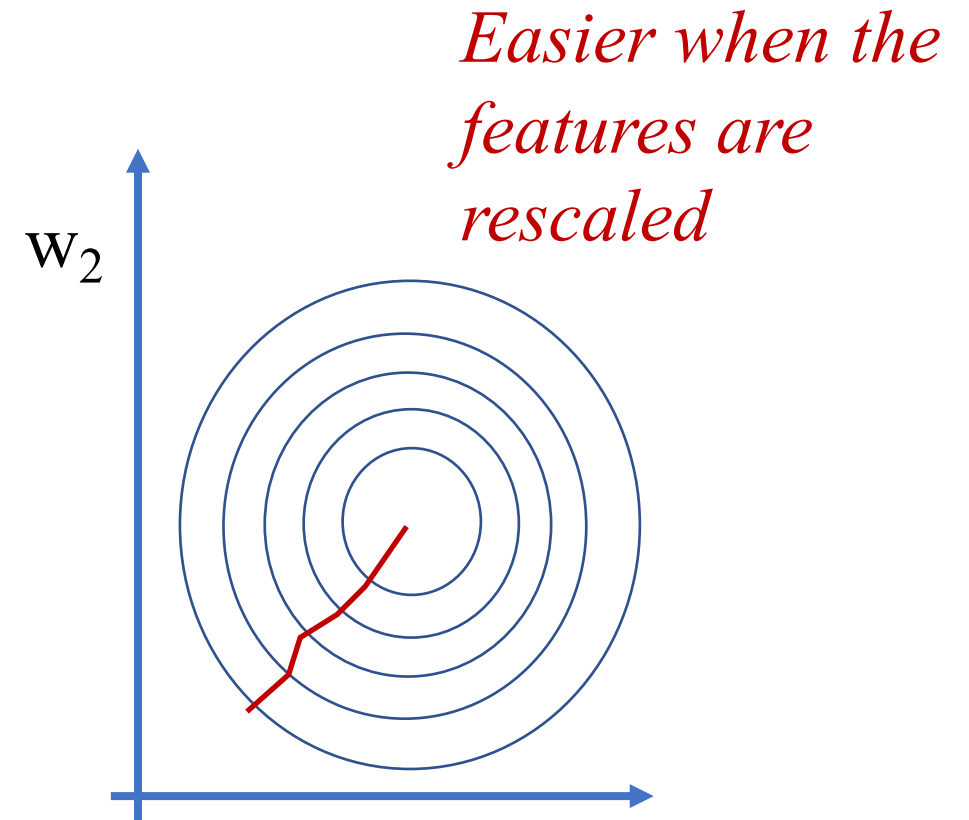
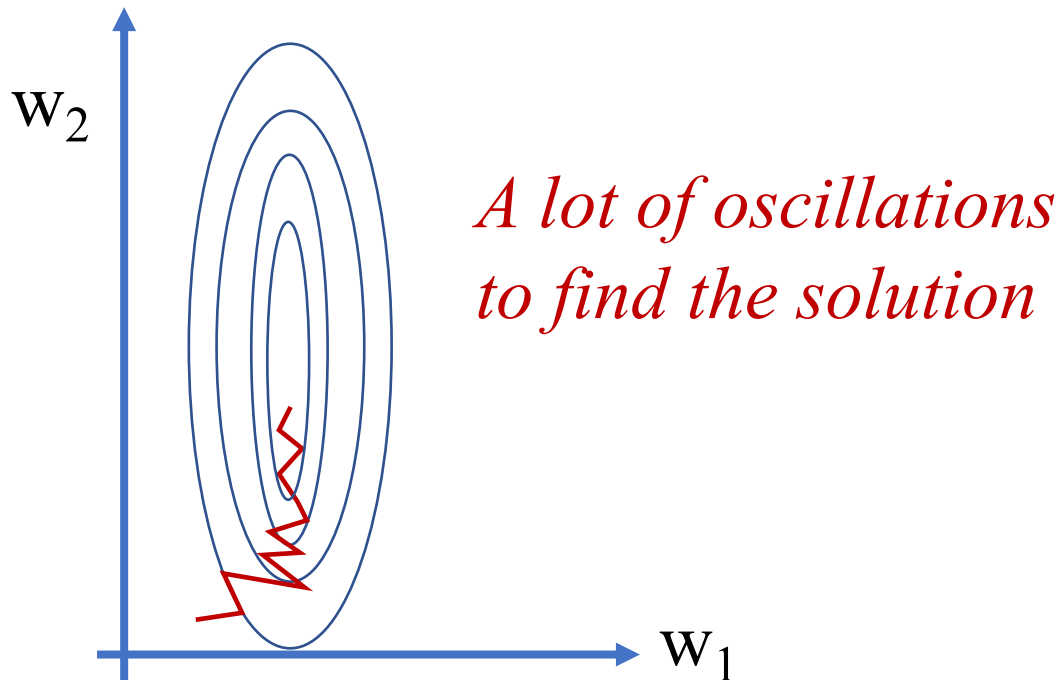
Practical Considerations

Feature Scaling

Need to scale features

72

- Gradient descent may not work as well, when feature values are not in the similar range of values.
- For example,
 - house area = 20-350m², but
 - number of bathrooms = 1-5



Rescaling also promotes interpretability

73

- If we don't rescale the features to have the same range (e.g., 0-1) then the weights for the linear regression don't mean much.

$$\text{Salary} = w_0 + 5\text{years_of_experience} + 2\text{number_of_friends} + \dots$$

- BUT, we all features have the same range, then the weights give you the **relative importance** of the features.

0-1 scaling

74

- Goal: to get every feature in the 0-1 range.
- Formula: $x_{\text{new}} = (x - \min(x)) / (\max(x) - \min(x))$

Mean normalization

75

- Another possible normalization is to remove the mean value of that feature and divide by some constant (either (max-min), or std).

$$x_new = (x - \text{mean}(x)) / (\text{max}(x) - \text{min}(x))$$

$$x_new = (x - \text{mean}(x)) / \text{std}(x)$$

Now it has zero mean and unit variance = **standardization**

Example: Assume x is the area of a house. Then we have:

$$x = (size - \text{mean}(size)) / (\text{max}(size) - \text{min}(size))$$

$$x = (size - 100) / (350 - 20) = (size - 100) / 330$$

Practical Considerations

Handling Categorical Variables
1-Hot encoding, Label encoding

Handling Categorical Variables

77

- So far we assumed the features were numerical.
- Many features will be categorical. How do we deal with those?

Handling Categorical Variables (Binary)

78

Categorical Binary (2 values) → convert into 0,1

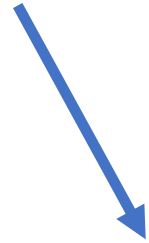


	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	PaymentMethod	Monthl
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	Electronic check	
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Mailed check	
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Mailed check	
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Bank transfer (automatic)	
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	Electronic check	
5	9305-CDSKC	Female	0	No	No	8	Yes	Yes	Fiber optic	Electronic check	
6	1452-KIOVK	Male	0	No	Yes	22	Yes	Yes	Fiber optic	Credit card (automatic)	

Handling Categorical Variables (>2 values)

79

Electronic check, Mailed Check,
Bank transfer, Credit card,...



	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	PaymentMethod	Monthl
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	Electronic check	
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Mailed check	
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Mailed check	
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Bank transfer (automatic)	
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	Electronic check	
5	9305-CDSKC	Female	0	No	No	8	Yes	Yes	Fiber optic	Electronic check	
6	1452-KIOVK	Male	0	No	Yes	22	Yes	Yes	Fiber optic	Credit card (automatic)	

Handling Categorical Variables (>2 values)

1-hot encoding

80

Electronic check, Mailed Check,
Bank transfer, Credit card,...



PAY_ECHECK	PAY_MCHECK	PAY_BANK	PAY_CREDIT	PaymentMethod	Month1
1	0	0	0	Electronic check	
0	1	0	0	Mailed check	
0	1	0	0	Mailed check	
0	0	1	0	Bank transfer (automatic)	
1	0	0	0	Electronic check	
1	0	0	0	Electronic check	
0	0	0	1	Credit card (automatic)	

Handling Categorical Variables (>2 values)

Label encoding

81

Electronic check, Mailed Check,
Bank transfer, Credit card,...



Electronic check: 1
Mailed check: 2
Bank transfer: 3
Credit card: 4
...

PAYMENT_METHOD	PaymentMethod	Month1
1	Electronic check	
2	Mailed check	
2	Mailed check	
3	Bank transfer (automatic)	
1	Electronic check	
1	Electronic check	
4	Credit card (automatic)	

1-Hot encoding vs Label encoding

Rule of thumb:

- When there are 10-15 distinct categorical values use 1-hot-encoding
- When there are more distinct categorical values use label encoding.
- In practice, you try both and see which one gives lower **test error. You use that one!**
- See also here an example:
 - <https://www.kaggle.com/alexisbcook/categorical-variables>

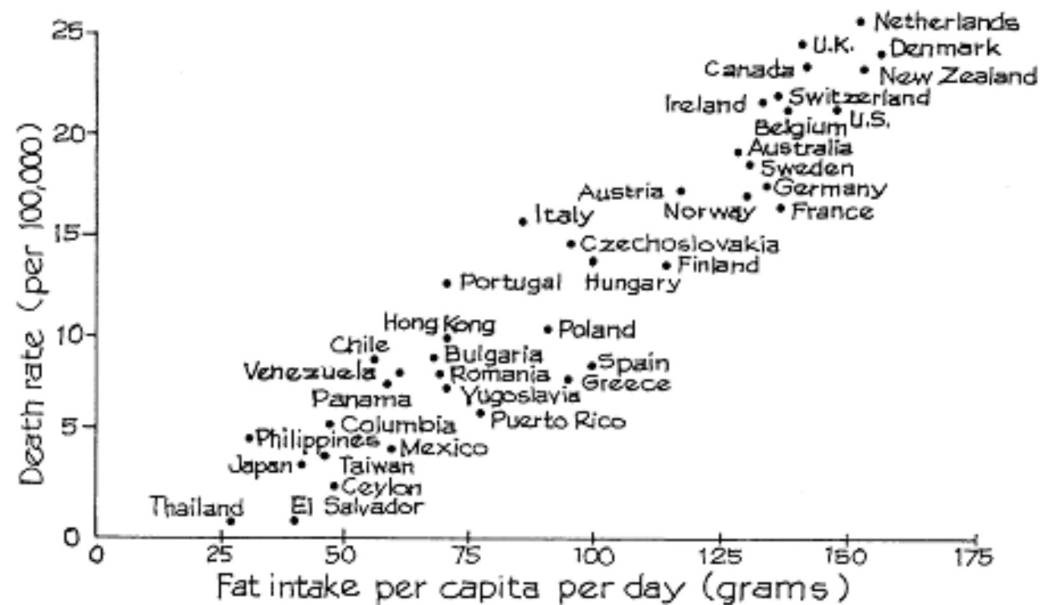
Correlation and Causation

Correlation/Association vs Causation

84

- When we find that some variables are associated, we should not fall into the trap of misinterpreting it as **causation**.

Figure 8. Cancer rates plotted against fat in the diet, for a sample of countries.

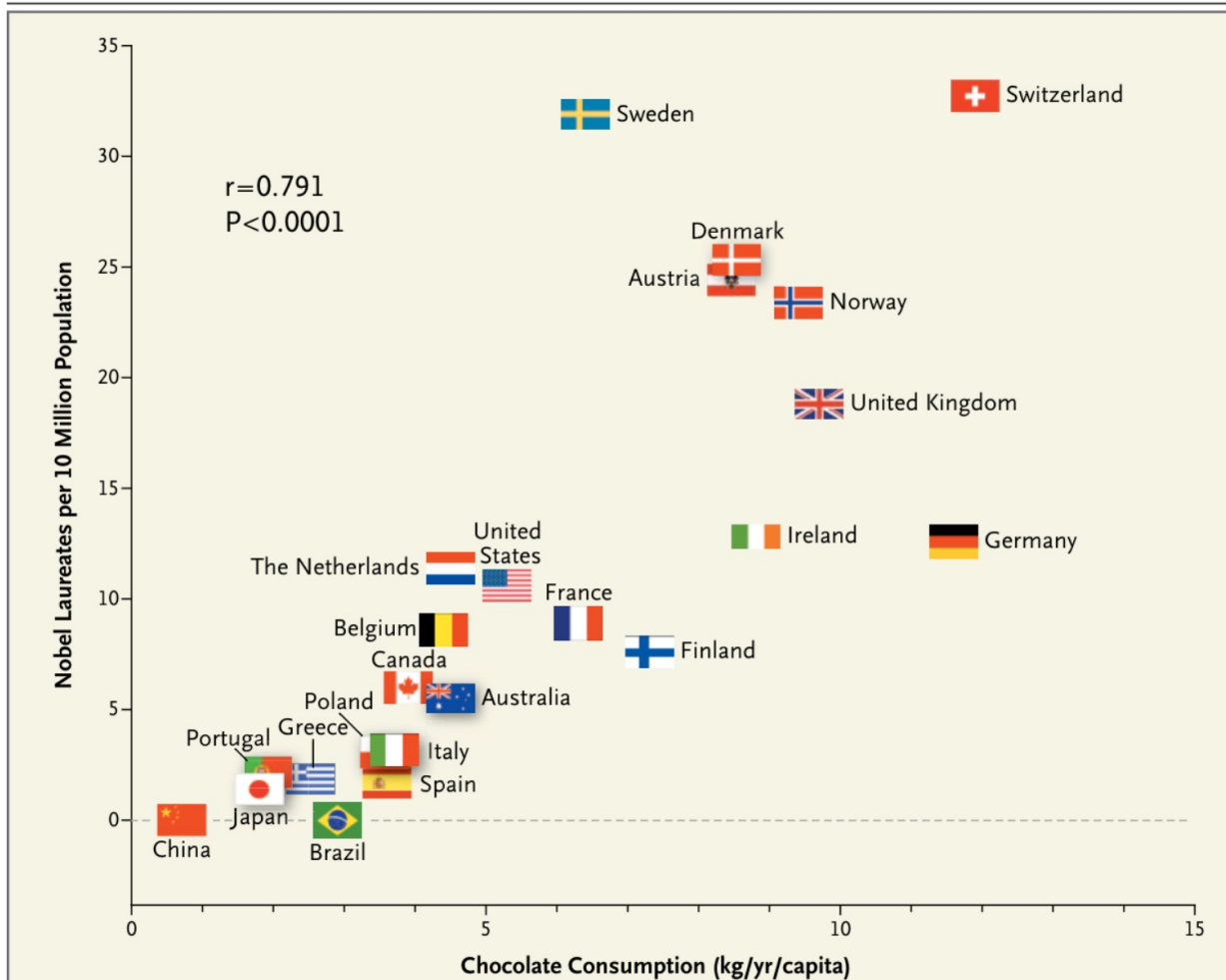


Does fat intake cause cancer?
This graph seems to provide such evidence.

Source: K. Carroll, "Experimental evidence of dietary factors and hormone-dependent cancers," *Cancer Research* vol. 35 (1975) p. 3379. Copyright by *Cancer Research*. Reproduced by permission.

Correlation/Association vs Causation

85



High-correlation between chocolate consumption and Nobel Laureates in a country...Switzerland rocks!

Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.





Correlation and Causation

Regression suggests that there is a correlation (especially for large R^2 values)

- **Correlation** – Values track each other
 - Height and Shoe Size
 - Grades and SAT Scores
- **Causation** – One value directly influences another
 - Education Level → Starting Salary
 - Smoking → Cancer

Correlation and Causation

Correlation does not imply causation

- Correlation can be result of causation from a hidden “confounding variable”
- A and B are correlated because there’s a hidden C such that $C \rightarrow A$ and $C \rightarrow B$

❖ Homeless population and crime rate

Confounding variable: unemployment

❖ Forgetfulness and poor eyesight

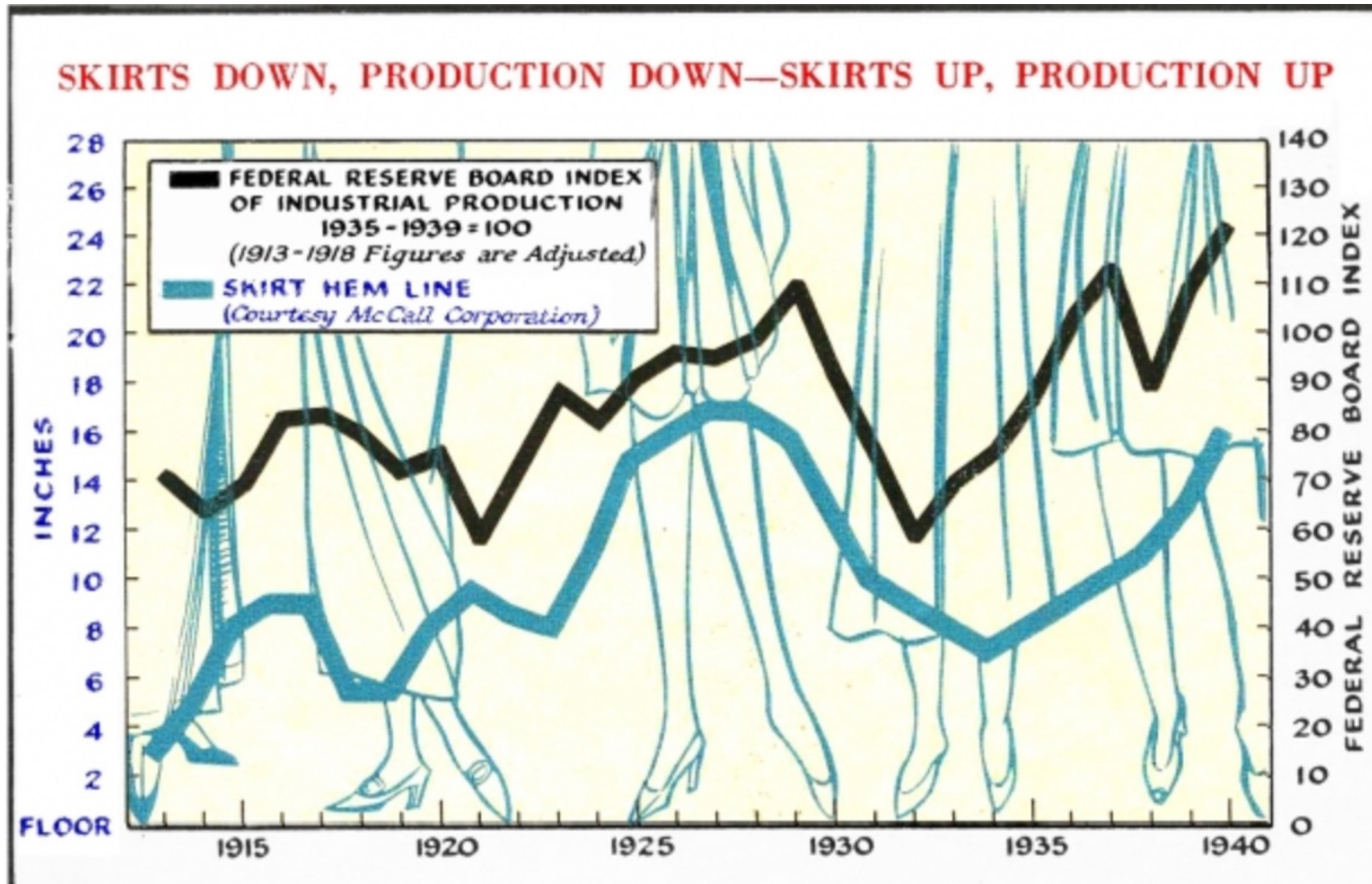
Confounding variable: age

Group Activity

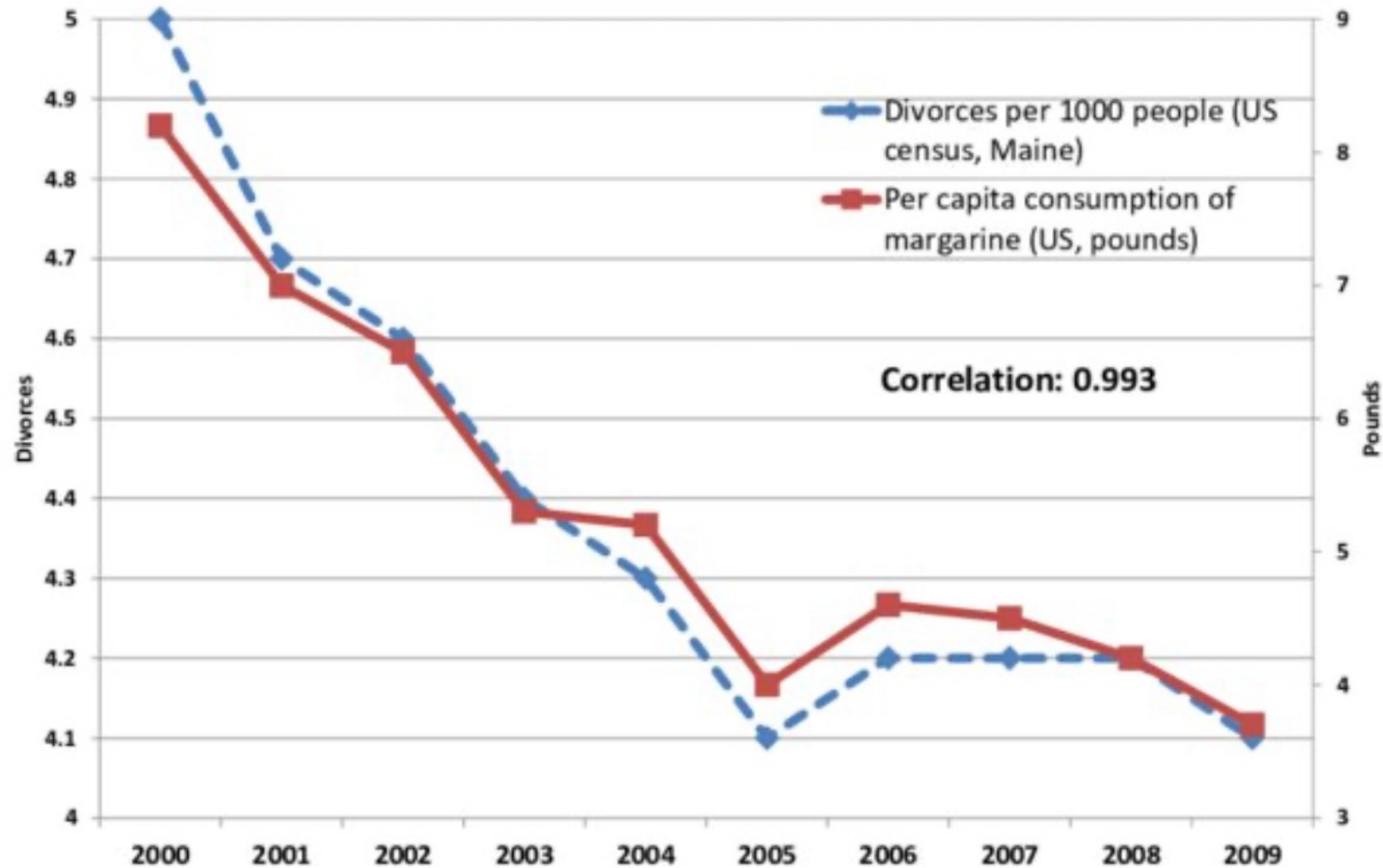
90

- Find some more cartoons that highlight the distinction between correlation and causation!

Surprising correlation #1



Surprising correlation #2



So remember. Correlation does not imply causation!

Terminology

Terminology

95

- Regression
- Linear regression
- Multi-linear regression
- Train(ing) error, test error
- Train/Test curves
- K-fold cross validation
- Rescaling
- 1-hot encoding
- Label encoding
- Causation vs correlation